# Image caption generation using Deep Q-learning framework

Dissertation Submitted In Partial Fulfillment Of The Requirements For The Degree Of

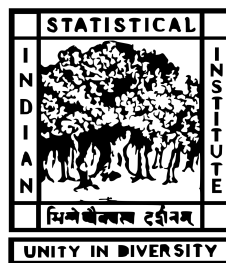Master of Technology
in
Computer Science

by

## Dipen Ganpat Rana
[ Roll No: CS1901 ]

Under the Guidance of

## Dr. Rajat K De
Professor
Machine Intelligence Unit(MIU)



### Indian Statistical Institute
### Kolkata-700108, India

# CERTIFICATE

This is to certify that the dissertation entitled **"Image caption generation using Deep Q-learning framework"** submitted by **Dipen Ganpat Rana** to Indian Statistical Institute, Kolkata, in partial fulfilment for the award of the degree of **Master of Technology in Computer Science** is a *bona fide* record of work carried out by him under my supervision and guidance. The dissertation has fulfilled all the requirements as per the regulations of this institute and, in my opinion, has reached the standard needed for submission.

**Dr. Rajat K De**
Professor,
Machine Intelligence Unit(MIU),
Indian Statistical Institute,
Kolkata-700108, India.

Dipen Rana,
MTech CS

# Acknowledgment

I would like to thank Dr. Rajat De, it was an absolutely great privilege and learning experience to work with him. He has been a constant source of support, starting from my first year. Discussions with him has always enlightened further path for me in dificult situations. I wish to express my sincere gratitude to all the research scholars in MIU Lab. Thanks to all the friends who have been there with me for the past two years. Lastly, I would like to thank my parents who have supported me and ensured my well being.

Dipen Rana,
CS1901
MTech CS
ISI Kolkata

# Contents

## Abstract

With the convolution neural network getting more and more popular in the last decade or so, thanks to the easy and affordable access to high computation power, it has been easy to process huge amount of data through these.

Machines are getting smarter and smarter with having capabilities in different domains. For example in Natural Language processing machines are now able to perform automatic language translations, voice question and answers systems, writing abstracts from the given text and many advanced processing. One such problem is caption generation for images. In this the machine has to generate a one line descirption for the input image. Various approaches tries to solve this problem with the use of LSTMs and attention networks.

In this project we propose a novel approach to solve the problem with the help of Deep Reinforcement learning. Our approach is based on the intution that in real life how the caption generation task is performed by an human language expert. A human expert will look at different parts of the image to get the keywords which should be present in the description of that image, and then it interprets the global information in the image to form the caption using the already sortlisted keywords. So hence we want to look at the local as well as global aspects in the image while formation of the caption, our model tries to mimic this approach where attention network and the LSTM network are used to capture local information and the reinforcement learning framework incorporates the global information in the process.

In this project we are trying to train a framework to correctly highlight the important parts of the image for classification task which can be used in caption generation task. The input image is feed into a CNN network and the features extracted from this will be given to the attention network which will decide what features to select based on attention scores. Basically, attention is a mechanism by which a network can weigh features by level of importance to a task, and use this weighting to help achieve the task for predicting the description for the image.

The dataset I will be using for this experiment is the MSCOCO dataset for the image caption generation task and training our model.

In the experiment, the input image is passed through a CNN to get the features, then it is passed through an attention network at each time step for getting the weighted attended features, which are used as the input with the predicted caption till now to the policy network to get the probabilities for selecting the next word. Based on the scores for the candidates which are calculated as reward for the Reinforcement learning is used to select he right word next. In this way more accurate predictions are generated since the LSTM is taking care for local information and the Reinforcement network taking care for the global information in the image.

# Chapter 1

# Introduction

The increasing computation and processing power and the ease of availability of these resources have led to the era of deep learning and Artificial inteligence. Now its so much easy to design and train much complex deep artificial networks with multiple layers and it is useful in solving many real life problems.

One such problem is the automatic image captioning. It refers to generating a textual description of an image by using some deep learning network. It combines both the image and textual processing to build a deep learning network for this task. Their are many potential applications of the image caption generation task in real life. For example, it can be used to search images related to a textual description or we can save the captions for images so it can be retrieved later based on the description. Other applications include such as recomendations in editing applications, usage in virtual assistants, for visually impaired persons, for social media and various other NLP applications.

With this project I tried to build a model which takes an image as input and it generates a one line description of that image as caption. There have been many approaches already developed for this task which uses encoder decoder model with attention networks and LSTM networks for predicting the captions.

The recent models have achieved great results for the image captioning task. It uses the encoder-decoder models which consist of a CNN network to get features representation of the image and the decoder part consist of LSTM network which will generate the words in sequence to form caption. The latest models are based on multi-layer Transformers [1]. Another types of models focus on different learning method. One such method is OSCAR [2], Object-Semantics Aligned Pre-training for Vision-Language Tasks, this method is based on the observation that the salient objects in the image can be accurately detected and are mentioned in the caption text. This model is briefly explained in the previous work section.

With this work we have tried to build a methodology which can work similar to human generated description so we can get similar results close to the naturally generated captions by some human expert. As human tries to write a description of an image, he mainly focuses on the different object which are their in the image and also what are the different actions performing in the image. Then based on these two things the human exprert tries to frame the sentence so that it is grammatically correct as well as describing the image in a nice detailed form.

Following this path, we formalized the method to generate the caption for an input image using the convolution neural networks, Attention networks [3], LSTMs with the encoder decoder networks

and a novel technique to train these networks by using the Deep Reinforcement Learning. Through this work, we have shown that how attention network along with the LSTM can be used to focus on different local objects in the image, and along side the task of capturing the global features and actions of the image can be captured using a novel decision making framework using Deep Q-Learning framework.

Also with this we have used the beam search technique so that if their is some error(grametical or syntactic) while framing the sentence then that can be correct in further steps. This also insures that the global aspects of the image is captured in the description sentence.

# Chapter 2

# Preliminaries and Previous Works

## 2.1 Backgound

Before introducing to my proposed framework, let us get some idea about the prerequisite concepts which are used in my solution.

### 2.1.1 Long Short Term Memory Networks

After RNN was introduced it was soon noticed that over time the gradient of the feed back signal would either vanish or explode. Schmidhuber et al. introduced Long Short Term Memory - Recurrent Neural Networks (LSTM-RNN) in 1997 [4] and improved it over time [5] to address the shortcomings of the regular RNN.

LSTM consists of special modules namely *Input Gate I()*, *Forget Gate F()* and *Output Gate O()*. Also, apart from hidden state, LSTM internally maintains the *Cell State* which helps it to keep track of the long term dependencies. Information can flow along the cell state unchanged, and if needed LSTM can easily add or remove information (gradient) with the help of input gate and forget gate respectively. And output gate helps the LSTM to generate the hidden state with the help of cell state.

For an input sequence $X = (x_1, x_2, ..., x_{n-1}, x_n)$, at time step $t$, the hidden state $h_t$ and cell state $c_t$ is calculated as shown in Eq. (2.1), where $f_t$, $i_t$ and $o_t$ are outputs from forget, input and output gates respectively, and $C()$ is an intermediary function with tanh activation.

$$
\begin{aligned}
f_t &= F(x_t, h_{t-1}) \\
i_t &= I(x_t, h_{t-1}) \\
o_t &= O(x_t, h_{t-1}) \\
c_t &= f_t \cdot c_{t-1} + i_t \cdot C(x_t, h_{t-1}) \\
h_t &= o_t \cdot H(c_t)
\end{aligned}
\tag{2.1}
$$

### 2.1.2  Gated Recurrent Unit

Motivated by the LSTM unit, in 2014 Cho et al. introduced Gated Recurrent Unit (GRU) which were much simpler to compute and implement [6]. Similar to what we saw in a LSTM unit, GRU also has gating units that modulate the flow of information inside the unit but with just the hidden state. No additional cell/memory state is present in GRU.

Apart from this, it also has fewer gates compared to LSTM, namely, *Reset Gate R*() and *Update Gate U*(). When reset gate is close to 0, the hidden state is forced to ignore the previous hidden state and reset with the current input only. This effectively allows the hidden state to drop any information that is found to be irrelevant later in the future, thus, allowing a more compact representation. On the other hand, the update gate controls how much information from the previous hidden state will carry over to the current hidden state .

Eq.(2.2) explains the calculation for the hidden unit $h_t$ at time $t$ for input sequence $X = (x_1, x_2, ..., x_{n-1}, x_n)$. $r_t$, $u_t$ are output of reset gate and update gate respectively and $H()$ is an intermediary function with tanh activation.

$$r_t = R(x_t, h_{t-1})$$
$$u_t = U(x_t, h_{t-1})$$
$$h_t = u_t \cdot h_{t-1} + (1 - u_t) \cdot H(r_t, h_{t-1}, x_t)$$

(2.2)

### 2.1.3  Encoder-Decoder Architecture

This architecture was introduced by Kalchbrenner et al. , Sutskever et al. , Cho et al. Encoder-Decoder (ED) [7] Architecture has gained lot of popularity in recent years and many of the state of the art models in Neural Machine Translation (NMT) [8] , Text Summarization , Image Captioning , etc. have ED Architecture at its core. Hereafter, we shall be referring to ED Architecture in context of the NMT.

Fundamental idea behind this architecture is, the encoder part reads the input (or a sequence of input) and condenses its meaning down to a fixed sized vector referred as context vector. This context vector is then fed to the decoder part which generates the desired results. The encoder part and the decoder part are generally recurrent neural networks and are both jointly trained in-order to maximize the probability of translation given a source sentence.

To understand this architecture more clearly, consider the following example. Let sample sentence $X$ be the source of length $n$ where $x_i$ is the word at $i^{th}$ position. Similarly, let $Y$ be the corresponding target (translation of $X$) sentence of length $m$ and $y_j$ be the $j^{th}$ word in $y$. Now, the encoder generates the hidden state $he_i = Encoder(x_i, he_{i-1})$ using $i^{th}$ word in $x$, i.e. $x_i$, and the hidden state from $i - 1^{th}$ step, i.e. $he_{i-1}$. This process is continued till $he_n$ is obtained, which is then used to initialize the hidden state decoder, i.e. $hd_0 = he_n$. It should be noted that $he_n$ or $hd_0$ is also called the context vector. Now the decoder takes in hidden state as $hd_0$ and a *special token* $< eos >$ as input and predicts $\hat{y}_1$ and outputs the next hidden state $hd_1$. This process continues in auto-regressive fashion, which can be summarized by the equation $\hat{y}_j, hd_j = Decoder(\hat{y}_{j-1}, hd_{j-1})$. The predictions by decoder continues and once $< eos >$ is obtained in the prediction, the decoder

4

stops and it symbolizes the end of prediction for the source sentence $x$. Figure (2.1) shows the high level overview of the architecture.
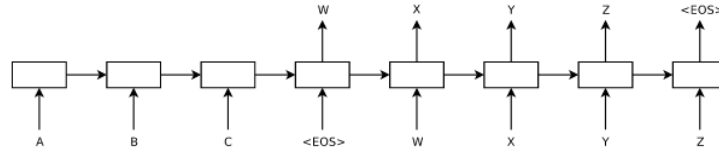


Figure 2.1: Model based on Encoder-Decoder Architecture for translating input sentence "ABC" and producing "WXYZ" as the output [?].

As one would notice, however long the input sample is, the context vector summaries all the information from it and is solely responsible for passing that information to the decoder for the final translation, due to which context vector itself became a bottleneck along with the long term dependencies between the source and the target sequence.

### 2.1.4 Attention network

With a neural network, it is considered to be a trial to mimic human brain actions in a simplified manner. The attention mechanism [9] is also such an attempt where the main focus is on few relevant things while ignoring the rest in a deep neural network.

In 2015, Bahdanau [10] came up with a simple but yet elegant idea where they proposed that in an encoder-decoder network, we can take the relative importance of each input word that can be taken into consideration for the Neural machine translation or NLP tasks. It was done using the attention mechanism in the encoder decoder architecture.



Figure 2.2: Attention mechanism in the encoder decoder architecture

The Figure (2.2) [11]shows diagram of the attention model from Bahdanau's paper. Here the LSTM generates a sequence of hidden state vectors i.e. h1,h2,h3...hT for each input sequence. In simple terms, all the vectors h1,h2....hT are hidden state representation of the T number of words in input. In simple encoder-decoder model, only the last hidden state representation is passed to the decoder module. Whereas, in case of attention network, the context vector $c_i$ for the decoder at time step t is generated using the weighted sum of the hidden state representations.

$$c_i = \sum_{j=1}^{T} \alpha_{ij} h_j \qquad (2.3)$$

here $\alpha_{ij}$ are the attention weights and these are computed using softmax function along with the other parameters of encoder-decoder model through backpropagation through time.

### 2.1.5  Deep Q-learning

Reinforcement learning [12] is very interesting idea in the artificial intelligence. It involves an agent, a set of states $S$ ,a set of actions $A$ and a reward associated with each combination of state and action. The objective of the reinforcement learning is to make optimal decisions using the past experience. The agent learns an optimal policy which helps in deciding how to take actions from a particular state so that the reward in response to the action taken in environment is maximized.

Q-learning [13] is a model-free and value based learning algorithm. Model free since it does not require a model of the dynamic environment to learn. This algorithm learns the values of an action in a particular state.

This Q-learning algorithm uses a data structure known as Q-Table.It is used to calculate the maximum expected future rewards for action at each state. This table is used by agent to select best action at each state. The values is Q-Table are filled using the Q learning algorithm.

Q-Learning uses a Q-function to eastimate the best policy for the task. The Q function is derived from the Bellman equation with the two inputs, here Q value for a given state s and action a at time step t is given below. The right side of the equation is the expected discounted cumulative reward given the state and action as s and a. Gamma is the discount rate.

$$Q^{\pi}(s_t, a_t) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + ... | s_t, a_t] \qquad (2.4)$$

The values of the Q-table are updated using the following set of equations, which is based on the bellman equation. The $\alpha$ is the learning rate.

$$Q'(s, a) = (1 - \alpha)Q[s, a] + \alpha Improvedestimate$$
$$Q'(s, a) = (1 - \alpha)Q[s, a] + \alpha(R + \gamma futurerewards)$$
$$(2.5)$$

$$futurerewards = Q[s', \arg\max_{a'} Q[s', a]]$$

The values of $\alpha$ and $\gamma$ are between 0 and 1. The high discount rate is like looking far ahead in future to estimate reward.

The Q-learning is simple yet quite powerful learning algorithm. It helps the agent to figure out the optimal choice of action at a particular state. But what if we have quite large number of states and actions, then following the vanilla Q-learning technique is not approprite to calculate the values of Q table. Instead we can use the machine learning model like neural network to approximate these Q values. This is the Deep Q-learning.

In Deep Q-Learning [14], we use a neural network to approximate the Q-value function. The input to the neural network is the state s, and the output is the predicted values for each action possible from this state. here the parameters of the neural network are learned through the backpropagation algorithm and the loss funciton is the mean squared error between the predicted and the target Q value.

## 2.2 Previous works

The problem of image captioning has been existing for a quite long ago. Their has been extensive work done in the past which has shown better results. Most of the proposed methods have used encoder-decoder kind of models, where the encoder is the convolution neural network CNN and the decoder consist of the recurrent neural network RNN. Here I have introduced couple of latest work done on this task.

### 2.2.1 Oscar: Object-Semantics Aligned Pre-training for Vision-Language Tasks

Large-scale pre-training methods of learning cross-modal representations on image-text pairs are becoming popular for vision-language tasks. While existing methods simply concatenate image region features and text features as input to the model to be pre-trained and use self attention to learn image-text semantic alignments in a brute force manner, in this method [2], we propose a new learning method Oscar1 , which uses object tags detected in images as anchor points to significantly ease the learning of alignments. Our method is motivated by the observation that the salient objects in an image can be accurately detected, and are often mentioned in the paired text. We pre-train an Oscar model on the public corpus of 6.5 million text-image pairs, and fine-tune it on downstream tasks, creating new state-of-the-arts on six well-established vision-language understanding and generation tasks.

In this study, we show that the learning of cross-modal representations can be significantly improved by introducing object tags detected in images as anchor points to ease the learning of semantic alignments between images and texts. We propose a new VLP method Oscar, where we define the training samples as triples, each consisting of a word sequence, a set of object tags, and a set of image region features. Our method is motivated by the observation that the salient objects in an image can be accurately detected by modern object detectors, and that these objects are often mentioned in the paired text. For example, on the MS COCO dataset [15], the percentages that an image and its paired text share at least 1, 2, 3 objects are 49.7model is pre-trained on a large-scale V+L dataset composed of 6.5 million pairs, and is fine-tuned and evaluated on seven V+L understanding and generation tasks.

### 2.2.2 Reflective Decoding Network for Image Captioning

State-of-the-art image captioning methods mostly focus on improving visual features, less attention has been paid to utilizing the inherent properties of language to boost captioning performance. In this paper, we show that vocabulary coherence between words and syntactic paradigm of sentences are also important to generate high-quality image caption. Following the conventional encoder-decoder framework, we propose the Reflective Decoding Network (RDN) [16] for image captioning, which enhances both the longsequence dependency and position perception of words in a caption decoder. Our model learns to collaboratively attend on both visual and textual features and meanwhile perceive each word's relative position in the sentence to maximize the information delivered in the generated caption. We evaluate the effectiveness of our RDN on the COCO image captioning datasets and achieve superior performance over the previous methods. Further experiments reveal that our approach is particularly advantageous for hard cases with complex scenes to describe by captions.

In this method, they propose the Reflective Decoding Network (RDN) [16] for image captioning, which mitigates the drawback of traditional caption decoder by enhancing its long sequential modeling ability. Different from previous methods which boost captioning performance by improving the visual attention mechanism , or by improving the encoder to supply more meaningful intermediate representation for the decoder , their RDN focuses directly on the target decoding side and jointly apply attention mechanism in both visual and textual domain.

Besides, we propose to model the positional information of each word within a caption in a supervised way to capture the syntactic structure of natural language. Another advantage in RDN is to visualize how the model inferences and makes word prediction based on the generated words. For instance, our RDN successfully decodes the word 'river' in Figure by referring to the previously generated words, especially the most relevant word 'bridge'.

The main contributions of this paper are four folds:

- They propose the RDN that effectively enhances the long sequential modeling ability of the traditional caption decoder for generating high-quality image captions.

- By considering long-term textual attention, we explicitly explore the coherence between words and visualize the word decision making process in text domain to show how we can interpret the principle and result of the framework from a novel perspective.

- We design a novel positional module to enable our RDN to perceive the relative position of each word in the whole caption and thereby better comprehend the syntactic paradigm of natural language.

- Our RDN achieves state-of-the-art performance on COCO captioning dataset and is particularly superior over existing methods in hard cases with complex scenes to describe by captions.

# Chapter 3

# Deep Reinforcement learning Methodology

In this paper we have tried to solve the image captioning task with some different approach. My approach is based on what actually the thought process is followed when some human is trying to give caption to an image. When a human expert is writing the caption for an image, it looks at the local as well as global aspects of the image and try to describe that in the caption. The thought process of a human while describing an image regarding the local and global aspects is as following.

**Local aspect** captures the local information from the image. A human expert first identifies different objects in the image by individually looking at certain parts of the image. These objects are identified as the keywords which must appear in the caption describing the image.

**Global aspect** focuses on the actions described in the image. It is critical for the formation of the caption with the correct alignment of the keywords and the filler words to complete a sentence. So, the human expert see the full image as frames the sentence with the keywords for different objects and the actions present in the image to form a good caption.

For example, here in the figure (3.1), the image shows the objects like a boy, baseball bat which can be seen as the local objects in the image. These are the keywords which a human can get from the local inference of the image. Also upon looking in the global picture of the image, the actions like swinging and the whole context of a baseball game can be inferred. So arranging all these words while having a global look of the image result in the caption "A boy swinging a baseball bat during a baseball game."

Through the following sections, We will describe how this intuition is carried out in my proposed Methodology for the image captioning task.

## 3.1 Proposed Methodology

In my proposed methodology, it is an encoder-decoder kind of methodology with the attention network and Q-learning. Here the attention network will highlight the important locations in the image thus capturing the local inferences form the image. On the other hand it also uses a scoring method to select the next best word for the caption from the dictionary which is trained using Q-learning technique. This handles the global inferences while generating a caption.

Figure 3.1: Sample image with the caption: A boy swinging a baseball bat during a baseball game.

The figure (3.1) shows the full methodology setup for my solution for the task. It is an encoder-decoder type of framework which consist of a CNN network as the encoder of the images into feature matrix which is feed to an attention network which give us the weighted encoding of the image feaures. This encoding is then passed through the decoder part which generates the caption word by word. It contains two different networks i.e. policy network and the value network.

Next I have explained each of the different modules and networks in the whole methodology.

### 3.1.1 Encoder CNN

The Convolution Neural Network use for image encoding can be any of the latest pretrained networks like ResNet, LeNet, VGG 16 etc. But for this task I have choose the VGG 16 network for
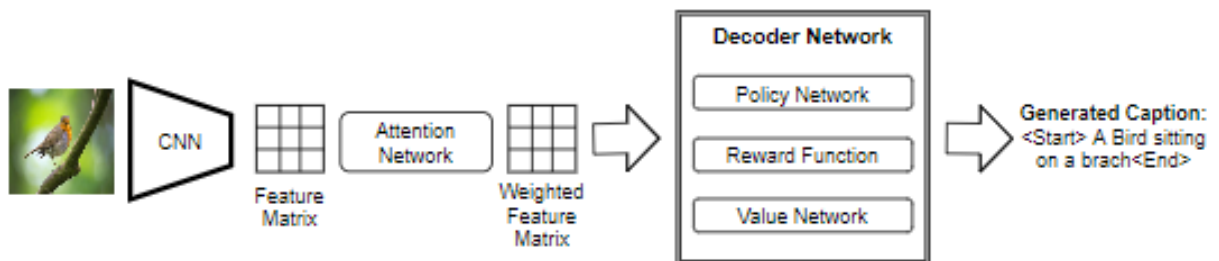


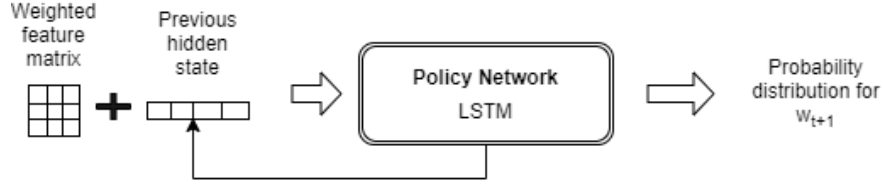Figure 3.2: Overview of the methodology for image caption generation

Figure 3.3: Overview of the policy network

encoding the images into features matrix. Since, we wanted the features which preserves the location information in the original image. Hence we take only the output from the last convolution layer in the VGG 16 network the output is taken from the 5th convolution-maxpool layer as the image embedding feature matrix.

The image encoding from the VGG-16 [17] is readily available for the images in the Microsoft Common Objects in context dataset MSCOCO. These features are then passed through the Attention network which highlights the important features from the feature matrix at each time step.

The attention parameters are the attention weights which are learned with the policy network parameters. It gives the feature matrix where some portions of the image features are highlighted by using the weights, these will be referred as weighted feature matirx which will be used in subsequent decoder networks.

### 3.1.2 Policy Network

The policy network is an LSTM network which is helps in selecting the best word next in the sequence of caption. The decoder part consist of the policy network, value network and the reward network which works together to generate the output caption for an image.

The main purpose of the policy network is to give us probability distribution for selecting the next word in the caption sequence. Before describing the policy network let me introduce some terminologies related to framing the task of image captioning into the decision making framework Reinforcement learning.

**State** will be referred as $s_t$ at time step t, it consist of the weighted encoding of image i.e. weighted feature matrix which we get from the attention network. This feature matrix is flatend and concatenated with the caption predicted until the current time step. For this the hidden state vector from the LSTM network (policy network) is taken. So state is combination of the image features and the predicted caption.

**Action** is selecting the next word from the dictionary, $a_t$. The policy network takes state as input and give us the probabilities for selecting an action $a_t$ which is the next word in the caption sequence. Figure (3.3) shows the policy network.

So the LSTM in the policy network gives us a probability distribution for all the words in the dictionary which are likely to be the next word in the sequence. The policy network is like an agent which learns an optimum policy to select the best word in the sequence given the input state. Later we will see that the top k number of words are taken into consideration for the beam

search algorithm with K number of candidate captions. This will be done by the value and rewards networks.

### 3.1.3  Reward Network

The Reward Network is used to compare the similarities between the image and the generated caption. I uses an encoding for the image as well as caption which is in same space and the similarity score is calculated based on these embedding. The reward network consist of linear layer which transform the feature matrix for an image which we get from CNN into a visual feature embedding, on the other hand to encode a sentence in the same format of embedding it uses a GRU network and that will give us the semantic embedding corresponding to the sentence caption.

For mapping the image into embedding space, we use the featrue vector $v$ from the CNN network used in the encoder, which is passed through a linear mapping layer denoted as $f_e$. For a sentence, its embedding features are taken as the last hidden state of the Gated Recurrent unit network, denoted as $h'_T(S)$. The following equations is denoting the computations.

$$VisualEmbedding, v_e = f_e(CNN_r(I))$$
$$SemanticEmbedding, s_e = h'_T(S)$$

(3.1)

$$Reward, r = \frac{v_e.s_e}{\|v_e\| \, \|s_e\|}$$

Here, $I$ and $S$ are the input image and the corresponding caption sentence. The rewards is calculated as dot product of normalized vectors $v_e$ and $s_e$. This is also called as the visual-semantic loss value.

The CNN is the same pretrained model as used in encoder. And linear mapping layer and the GRU are trained using on the same set of images and captions form COCO dataset.

The reward score is used only for training the value networks multpile layers. while training the value network, the visual-semantic loss [18] is used as the target value for the image-caption pair to train the MLP model using the backpropagation algorithm.

### 3.1.4  Value Network

Value network is responsible for evaluating the reward score for all possible extensions of the current state of the caption. It serves as the global guidance and lookahead by maintaining a set of candidates in a beam search. The value network consist of the MLP, which has only two layers and it takes the input as the concatenation of the image features and the partially generated caption. The output of the value network is a scalar score which compares the goodness of caption based on the current and future states.

From the Policy network we take some top words form the distribution for $w_{t+1}$ word in the caption. Then the value network evaluates the scores for each combination of candidate captions and it will

maintain the decreasing order of candidate captions in the beam. At last, it will output the caption with the best score in the beam as the generated caption. This ensures that lookahead information is used while framing the sentence. This ensures that even if their is a bad word selection at some timestep t from policy network, their is still chance to correct this mistake in subsequent timesteps as we are maintaining k number of candidate captions and the other best caption will have a higher score and moved up in the beam.

The training of the Value network is based on the Q-learning which is based on the visual-semantic loss as the reward function. The MLP in value network is trained using backpropagation over the Mean-Squared error loss between the value given by the value network and the reward score from the reward network. This value is minimized using the Adam optimizer.

## 3.2   Experiment

This section explains the details of the experiment performed with the proposed architecture. All the code are run on the Google Colab server with 1 CPU and 1 GPU. First we will discuss the dataset used for this experiment and the implementation details and then we will compare the results with the state-of-the-art models for image captioning.

### 3.2.1   Dataset

As it is already mentioned that dataset used is from the MSCOCO dataset [15]. We have taken a subset of this data. So our training set is having around 50,000 images-caption pairs and the validation set contains around 2000 pairs. We have taken only a subset of the whole dataset because of the avalilabity of limited computaion power on google colab with the restrictions on time for using GPU power on the free account. So with the reduced set, I was able to train all the models in my methodology and the parameters are saved in PT files for further use so it not require to train model everytime using it.

This has also impacted the result of this experiment, as less number of captions for training means that the dictionary size will be small and their are chances of encountering unknown words in the validation time. But still the results are quite satisfactory to prove that the methodology is capable of producing comparable results given high computaion power.

### 3.2.2   Training

The training of the whole framework is done in a step-wise manner with taking into consideration the dependecies of different models in the architecture. Here I have explained the training of each of these different models.

First the **encoder CNN** is taken as the VGG-16. Fortunately the pretrained model of VGG-16 is avilable for the MSCOCO dataset. Even the fearues from the last convolution layer of the VGG-16 architecture are available with the dataset. For this experiment I have used the same features with PCA applied on 4096 dimensional features which are reduced to 512.

The **Attention network** and the **Policy network** are trained together as these networks are closely coupled with each other. The tainable model parameters here are the attention weights (*Alphas*), and the LSTM parameters U,V,W and the gates. The training algorithm used is the

Backpropagation through, and optimzer used is the Adam optimizer and loss funciton as cross entropy loss.

**Reward network** is trained independently which is used for training the Value network. The parameters here are the linear layer weights for visual embedding and the GRU parameters. Optimizer is the Adam optimizer and the loss function is the visual semantic loss.

For training the **Value network**, we require the output from the Policy network and reward score form reward network. Adam optimizer is used with the MSE between the Value score and the reward. The Q-learning technique is used here to train this model.

Once all these models are trained, the best model parameters are stored in the Kodac Precision Transform(.pt) files for evalutaion and further use.

### 3.2.3    Evaluation metrics

For evalutaion of the proprosed architecture I have used all the standard scores which are BLEU scores, METEOR, ROUGE-L, CIDEr. out of these I have focused on the BLEU scores, Meteor and CIDEr scores which are briefly defined below.

**Bilingual Evaluation Understudy** (BLEU) [19] score is quality metric for machine translations based on the string-matching algorithm. It attempts to measure the correspondence between a machine translation and the human translation. It takes one or more human reference translations and compare it against the machine translation. Differnet bleu scores like 1,2,3 etc are just the n-gram matching of a specific order between reference and generated outputs.

**Metric for Evaluation of Translation with Explicit ORdering** (METEOR) [20] is a metric based on the harmonic mean of the unigram precision and recall, with recall weighted higher than precision. It also take care for the stemming and synonyms matching in the reference and generated sentences.

**Consensus-based Image Description Evaluation** (CIDEr) [21] it measures the similarity of a generated sentence against a set of ground truth sentences written by humans. It uses the tf-idf metric to aggregate statistic for n gram across the data. It means that words present across many captions is less informative, thus less weight is given to them in evaluating the similarity. with using similarity, the notion of grammatical correcteness, saliency, importance and accurracy are inherently captured by this metric.

## 3.3    Results

The Table (3.1) is provides a summary of the results of my model for the image captioning task. The table shows the results for 2043 images from the test set.

The scores given by my framework are not very high, but still proves the effectivness of the framework in generating good captions for some of the cases. It is evident that rigorous training and parameter tuning can bear quality results from this proposed architecture. The Figure shows some of the captions generated for the test image samples by using this framework.

In Table (3.2), I have compared the results from my model with the state-of-the-art models for

|        | Test Set |
|--------|----------|
| BLEU-1 | 39.97 |
| BLEU-2 | 26.04 |
| BLEU-3 | 19.30 |
| BLEU-4 | 15.61 |
| METEOR | 21.87 |
| ROUGE_L | 38.44 |
| CIDEr | 64.72 |

Table 3.1: Evaluation scores for the test set of MSCOCO dataset with our deep reinforcment learning framework.

image captioning task. We have not got comparable results against these models but still, the methodology has shown good signs of capturing local and global details for some of the samples. If we consider the table, we can see that BLEU 4 score for the OSCAR and other models is around 36.5, and also this score for other models like RDN and Google NIC is over 30, but here our model is showing score of 15.61, this is because since our model is trained on less data, it not performing good for captions having large lengths, so the BLEU 4 score which is based on 4-grams statistic is very low. Also the other scores like METEOR for other models is in range of 25 to 30, our model is also showing little comparable results here with the score of 21.87. The CIDEr for our model is 64.72 whereas the most successful model having this score of 127.8 .

|                   | Bleu-1 | Bleu-2 | Bleu-3 | Bleu-4 | METEOR | Rouge-L | CIDEr |
|-------------------|--------|--------|--------|--------|--------|---------|-------|
| Google NIC [22]   | 71.3   | 54.2   | 40.7   | 30.9   | 25.4   | 53.0    | 94.3  |
| RDN               | 77.5   | 61.8   | 47.9   | 36.8   | 27.2   | 56.8    | 115.3 |
| OSCAR             | -      | -      | -      | 36.5   | 30.7   | -       | 127.8 |
| Ours              | 39.97  | 26.04  | 19.30  | 15.61  | 21.87  | 38.44   | 64.72 |

Table 3.2: Comparision of results from our model with the latest state-of-the-art models.

The Figure (3.2) shows some of the good sample results generated from the proposed methodology. In these samples the generated captions are very much related to the images and also very much similar to the actual captions. It may be noted that the evaluation metrics used for evaluating this kind of task is based on the reference captions, and it may be posible to write a caption for an image in a different style with reference caption which will give a very low score, although the generated caption is describing the image very nicely.

Actual:<START> a man <UNK> tennis <UNK> to waiting tennis players <END>
Generated: <START> a close shot of a female tennis player holding a tennis racket <END>

Actual:<START> young boys are playing <UNK> on a dirt field <END>
Generated: <START> a boy swinging a baseball bat during a baseball game <END>

Actual:<START> an elephant in the water inside its <UNK> <END>
Generated: <START> two elephants standing in the water together <END>

Actual:<START> three street signs sit atop a one way sign <END>
Generated: <START> there is a street sign next to each other <END>

Actual:<START> a young man standing on a tennis court holding a racquet <END>
Generated: <START> a man is on a court with a tennis racket <END>

Actual:<START> a brown and white cat <UNK> in front of a glass vase containing a plant <END>
Generated: <START> a black and white cat is <UNK> next to a vase <END>

Figure 3.4: Some sample captions generated by the proposed framework form the testset.

16

# Chapter 4

# Conclusion and Future Works

Through this work I have tried to purpose a novel architecture for the image captioning task. The intution is to mimic the process for writing captions in real life by using machine learning and reinforcement learning framework. The proposed framework is based on the encoder decoder kind of architecture with attention network and the Deep Q learning technique. We defined a decision making framework and formulated the task of generating image related sentences into Reinforcement learning framework. I have defined the policy network, reward network and value network which work together to generate good quality caption. In this framework the attention network and the policy network processes the local information from the image and the value network along with beam search technique incorporated the global information from the image into the caption.

The results I got from this methodology is good enough to show that the method followed here is somewhat benifical and this can lead to better comparable results if rigorous training and tuning is done. But the sample results shown in the previous section shows that the generated captions are sometimes way better than the actual captions. This conclude that this idea is indeed interesting to work upon and future research can be done in this direction to generate better results for the image captioning task. Also similar kind of idea can be applied to other areas in NLP such as visual question answering system, scene-graph generation or textual commentary generation over video.

In Future, the architecture can be improve by using the better value estimator for the value network. Also by training the architecture on a bigger dataset will improve the predicted captions quality.

# Chapter 5

# Bibliography

[1] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I. Attention is all you need., NeurIPS 2017.

[2] Chunyuan Li Pengchuan Zhang Xiaowei Hu Lei Zhang Lijuan Wang Houdong Hu Li Dong Furu Wei Yejin Choi Xiujun Li, Xi Yin and University of Washington Jianfeng Gao, Microsoft Corporation. Oscar: Object-semantics aligned pre-training for vision-language tasks, 2020.

[3] Marco Lippi Andrea Galassi and Paolo Torroni. Attention in natural language processing, 2019.

[4] Jurgen Schmidhuber Sepp Hochreiter. Long short-term memory, 1997.

[5] Franc̦oise Beaufays Hașim Sak, Andrew Senior. Long short-term memory recurrent neural network architectures for large scale acoustic modeling. 2014.

[6] Bart; Gulcehre Caglar; Bahdanau Dzmitry; Bougares Fethi; Schwenk Holger; Bengio Yoshua Cho, Kyunghyun; van Merrienboer. "learning phrase representations using rnn encoder-decoder for statistical machine translation", 2014.

[7] Nal Kalchbrenner Lasse Espeholt Karen Simonyan Aaron van den Oord Alex Graves Koray Kavukcuoglu. Neural machine translation in linear time, 2016.

[8] Felix Stahlberg. Neural machine translation: A review and survey, 2019.

[9] Ryan Kiros Kyunghyun Cho Aaron Courville Ruslan Salakhutdinov Richard Zemel Yoshua Bengio Kelvin Xu, Jimmy Ba. Show, attend and tell: Neural image caption generation with visual attention, 2015.

[10] Yoshua Bengio Dzmitry Bahdanau, Kyunghyun Cho. Neural machine translation by jointly learning to align and translate. 2015.

[11] Introduction to bahdanau attention and luong attention attention mechanism.

[12] Andrew W. Moore Leslie Pack Kaelbling, Michael L. Littman. Reinforcement learning: A survey, 1996.

[13] Beakcheol Jang, Myeonghwi Kim, Gaspard Harerimana, and Jong Wook Kim. Q-learning algorithms: A comprehensive classification and applications. *IEEE Access*, 7:133653–133667, 2019.

[14] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013.

[15] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015.

[16] Lei Ke, Wenjie Pei, Ruiyu Li, Xiaoyong Shen, and Yu-Wing Tai. Reflective decoding network for image captioning, 2019.

[17] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015.

[18] Jon Shlens Samy Bengio Jeff Dean Marc'Aurelio Ranzato Tomas Mikolov Andrea Frome, Greg S. Corrado. Devise: A deep visual-semantic embedding model, 2013.

[19] Todd Ward Kishore Papineni, Salim Roukos and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation, 2002.

[20] Abhaya Agarwal Alon Lavie. Meteor: An automatic metric for mt evaluation with high levels of correlation with human judgments, 2007.

[21] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation, 2015.

[22] S. Bengio O. Vinyals, T. Alexander and D. Erhan. Show and tell: a neural image caption generator, 2014.