

Article

Dual-Modal Transformer with Enhanced Inter- and Intra-Modality Interactions for Image Captioning

Deepika Kumar ^{1,*} , Varun Srivastava ², Daniela Elena Popescu ³  and Jude D. Hemanth ⁴ 

¹ Department of Computer Science & Engineering, Bharati Vidyapeeth's College of Engineering, New Delhi 110063, India

² Department of Computer Science and Engineering, Thapar Institute of Engineering and Technology, Patiala 147004, India; varun.srivastava@thapar.edu

³ Faculty of Electrical Engineering and Information Technology, University of Oradea, 410087 Oradea, Romania; depopescu@uoradea.ro

⁴ Department of Electronics & Communication Engineering, Karunya Institute of Technology and Sciences, Coimbatore 641114, India; judehemanth@karunya.edu

* Correspondence: deepika.kumar@bharativedyapeeth.edu

Abstract: Image captioning is oriented towards describing an image with the best possible use of words that can provide a semantic, relatable meaning of the scenario inscribed. Different models can be used to accomplish this arduous task depending on the context and requirement of what needs to be achieved. An encoder–decoder model which uses the image feature vectors as an input to the encoder is often marked as one of the appropriate models to accomplish the captioning process. In the proposed work, a dual-modal transformer has been used which captures the intra- and inter-model interactions in a simultaneous manner within an attention block. The transformer architecture is quantitatively evaluated on a publicly available Microsoft Common Objects in Context (MS COCO) dataset yielding a Bilingual Evaluation Understudy (BLEU)-4 Score of 85.01. The efficacy of the model is evaluated on Flickr 8k, Flickr 30k datasets and MS COCO datasets and results for the same is compared and analysed with the state-of-the-art methods. The results shows that the proposed model outperformed when compared with conventional models, such as the encoder–decoder model and attention model.

Keywords: attention model; encoder–decoder model; multi-modal transformer; BLEU score; beam search



Citation: Kumar, D.; Srivastava, V.; Popescu, D.E.; Hemanth, J.D. Dual-Modal Transformer with Enhanced Inter- and Intra-Modality Interactions for Image Captioning. *Appl. Sci.* **2022**, *12*, 6733. <https://doi.org/10.3390/app12136733>

Academic Editors: Mourad Oussalah and Rachid Jennane

Received: 28 April 2022

Accepted: 25 June 2022

Published: 2 July 2022

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Image captioning has proved its advantage at the naive level by being a virtual assistant or a means of understanding for visually challenged people [1,2]. The task of image captioning has turned out to be difficult due to varied interpretations of the image by different individuals; thus, making it complex enough to come out with an accurate description of the image with the best possible use of words, and this is even a definitive way of caption generation for images [3]. The model's output is envisioned to describe what is shown in the image utilising the object's properties, but replicating this behaviour in an artificial system is a tedious and lengthy task; hence the incorporation of intricate techniques, such as deep learning, to solve the task comes into application [4]. Image captioning as compared to image recognition is a complex process, as the former allows the machine to present a sentence in the correct manner encompassing the relevant details of the scene [5]. Advancement in the process of image captioning leads to a better development or usage of a model to achieve significant results and analysis. When the objects present in the considered image are detected and their relationship has been identified, then a caption is generated for the image. The sequence of the words matters the most from which the captions are generated since a readable and grammatically correct sentence is important for understanding [6].

To perform a comparative analysis and concise review of different image captioning methodologies, a detailed study and implementation of both attention-based and non-attention-based image captioning methodologies had been performed. The BLEU score or Bilingual Evaluation Understudy has been taken as the evaluation criteria, which is one of the most trusted evaluation metrics for the quality of text generated by the machine [7,8]. Flickr 30k dataset has also been used on the above architecture that includes about 30,000 images with each image used for five different captions leading to a total of 158k captions providing a benchmark in the domain of sentence-based image captioning [9]. A simple attention-based model, as previously mentioned, can be gruelling at times purposely when the objective laid has abstract word identification in it, which leads to a semantic gap between the words and the vision. To bridge this gap, the transformer model has come into extrication, which can attain semantic and scenic information simultaneously [10]. Hence, an approach for image captioning based on dual-model transformer architecture is proposed, with iterative answer prediction using a dynamic pointer network (DPN) [11]. Given an image, n -dimensional feature vectors are extracted from two modalities, i.e., the text and object. These feature vectors are then passed as an input to a multi-layer transformer which outputs a reinforced representation of these vectors in a common semantic space through intra and inter interactions among the two modalities [12]. Thus, a caption is then generated word-by-word in an auto-regressive manner through the transformer using the dynamic pointer approach.

The proposed work concatenates two types of embeddings: one is generated using the image and another is word embedding of that image to feed a given encoder. This dual mode approach creates a highly extensive feature vector for the training of the encoders. The major novelty lies in the extraction of feature vectors fed into the encoder. To form the encoding from the image, two embeddings, one based on the objects in the image and another by using an Inception-V3 model are concatenated and a final image-based embedding is generated. Further, a Faster Regions with CNN (FRCNN) has been used along with a region prediction network (RPN) to generate object-based encoding from the images. This generates a detailed feature vector. The use of a highly extensive framework for feature extraction is responsible for a detailed and in-depth caption generation along with good BLEU, METEOR and ROUGE scores when compared with related techniques. Furthermore, the BLEU scores have been computed by using three different algorithms, i.e., beam search (with $K = 3$ and 5) and greedy algorithm.

This paper is organised in the following manner. Section 2 presents the related work and background for the proposed work; Section 3 describes the step-by-step methodology used; and Section 4 discusses the results obtained followed by the conclusion.

2. Background

In one of the previous work studies, the image captioning task had been carried out using the Flickr 30k dataset and the captioning system had been implemented using the long short-term memory (LSTM) network which sought to achieve a median rank (Med -R) of 5. An attention-based model has been used that combines both bottom-up and top-down strategies. Recurrent neural network (RNN) has been used that gives Bilingual Evaluation Understudy Score (BLEU)-2 and BLEU-3 scores of 50.4 and 35.7, respectively. One of the studied works proposed the methodology of merging systems that primarily takes into consideration partial captions, which are subsequently merged with image vectors to generate the desired captions, which also incorporates the RNN hidden state vector [13].

Visualising or seeing an image does not demand much effort, whereas, on the other hand, describing the scene in the image requires a huge effort from the human side. There has been a potential increase in the use of Flickr 8k and Flickr 30k datasets, which comprise about 8000 images and 30,000 images, respectively, purely justifying their naming convention [14]. The Flickr 30k dataset also depicts its advantageous behaviour in the domain of text grounding tasks, performing relatively well as compared to other conventional datasets by approximately 3.088%. Flickr and MS COCO simulate explosive growth in

different deep learning tasks and their respective applications. Faster RCNN can also be incorporated to generate text descriptions about a specific image in the sequence of LSTM and RNN, which resulted in a BLEU-1 score of about 59.8 [15]. The Flickr 30k dataset has established its importance and advantage over the field of automatic image captioning (AIC), preferred to depict the most remarkable data out of images by finding the relationship among different objects present in the image [16]. The half and half bidirectional LSTM approach with CNN [17] has also depicted significant outcomes for picture extraction and captioning on Flickr complete datasets, thereby achieving a true positive ratio of 86% and a false positive ratio of 10%. Extending the image captioning from English to the Arabic language, the authors have proposed a deep recurrent neural network—long short-term memory RNN-LSTM model concatenated with CNN [18], which uses Google translated Arabic captions for training the model. Furthermore, despite having a small Arabic version of Flickr and MS COCO, the model can achieve a BLEU-1 score of 46.2. The authors improvised on the conventional transformer model and introduced the object relation to the transformer model. The proposed methodology shows a significant increase in results. The CIDE-rand BLUE-4 scores achieved were 128.3 and 38.6, respectively [19]. The authors in [20] proposed a fully attentive image captioning methodology. The proposed architecture uses a multi-layer encoder–decoder architecture for generating image regions and output sentences. The MS COCO dataset has been used for experimentation and achieved a BLUE-4 score of 39.7. Fei et al. [21] proposed an attention-aligned transformer model for image captioning; image mask operation has been added to improvise on results. The proposed methodology achieved a BLUE-4 score with a value of 39.8 and CIDEr value of 133.9. S. Yan et al. [22] proposed a hierarchical approach employing the GAN framework as well as RL optimisation. This three-step training strategy produces high-quality captions achieving a BLEU-1 score of 73.03. For conducting object detection in an image, an object detection model called Faster R-CNN is used, which is trained on the COCO 2017 dataset [23]. This model consists of two major modules, the primer and one convolutional layer for proposed regions. The other module works as the detector by using these proposed regions to generate region proposals. From the coordinates, a 4-dimensional location feature vector is extracted for the top number of detected objects, after grouping the boxes that correspond to the same location and having a minimum score threshold [24]. ChenLong et al. [25] proposed a dubbed CNN model called spatial and channels attentions in convolution neural network (SCA-CNN), which tunes sentence generation in feature maps at multiple levels. The architecture delivers a BLEU-1 score of 66.2 on the Flickr 30k dataset. The authors in [26] proposed a hierarchical attention network (HAN), which works synchronously on hierarchical features in a pyramid fashion and successfully outputs accurate captions indicated by a BLEU-1 score of 80.9 and a CIDEr score of 121.7.

Zhou et al. [27] proposed a text-conditional attention mechanism that takes its input from guiding long short-term memory (gLSTM) paired with CNN fine-tuning. The outcomes on the MS COCO dataset yielded a BLEU-1 score of 71.6. The authors in [28] proposed a novel algorithm motion CNN for removing motion features for generating image captions, which minimise the accuracy of image captions. It enhances the performance of caption generation and achieved a BLEU-1 score of 75.9. Quanzeng et al. [6] proposed a combined top-down and bottom-up approach that fuses semantic concepts into hidden states and outcomes of RNN. The method achieves a BLEU-1 score of 71.9 on the Flickr 30k dataset. The authors proposed a stand-alone convolutional neural network against the conventional LSTM–RNN pairing [29]. Defying the sequential nature attained a BLEU-1 score of 72.5 on the MS COCO dataset. Lu Jiasen et al. [30] put forward an adaptive attention model, which toggles between spatial attention and visual sentinel. The LSTM extension, when evaluated on the Flickr 30k dataset, yields a BLEU-1 score of 67.7.

There can be different captions for a particular image based on various human interpretations, which often makes the task of image captioning an ambiguous one. To resolve this ambiguity, beam search can be used to derive sampling for different captions [31,32]. Beam

search works upon a set of captions intending to find the most suitable caption having the highest posterior probability and this sole requirement is having a wide beamwidth. During the decoding process, the captions are generated for an image through the output sequences with beam search by setting a particular value of beam width [33].

3. Proposed Methodology

There are a significant number of convenient development libraries, the most popular being Pytorch [34] and Tensorflow [35] along with datasets largely labelled as Flickr and MS COCO that simulate an explosive growth in different deep learning tasks and their respective applications. In the proposed work, a dual-modal transformer has been incorporated to generate accurate image captions for the MS COCO dataset. The embeddings used are obtained by concatenating the embeddings from object detection, the Inception V3 model and the cleaned captions from the MS COCO dataset. To generate the embeddings, CNN and geometry features have been used. This is done to obtain an enhanced feature vector based on two modalities viz. image-based embedding and text-based embedding. Using CNNs to extract features for a multimodal transformer used in image captioning has been justified in [5], and thereby Inception V3 has been used in the proposed approach. The steps for extraction of multimodal features and the corresponding detailed encoding pedagogy are shown in Figure 1.

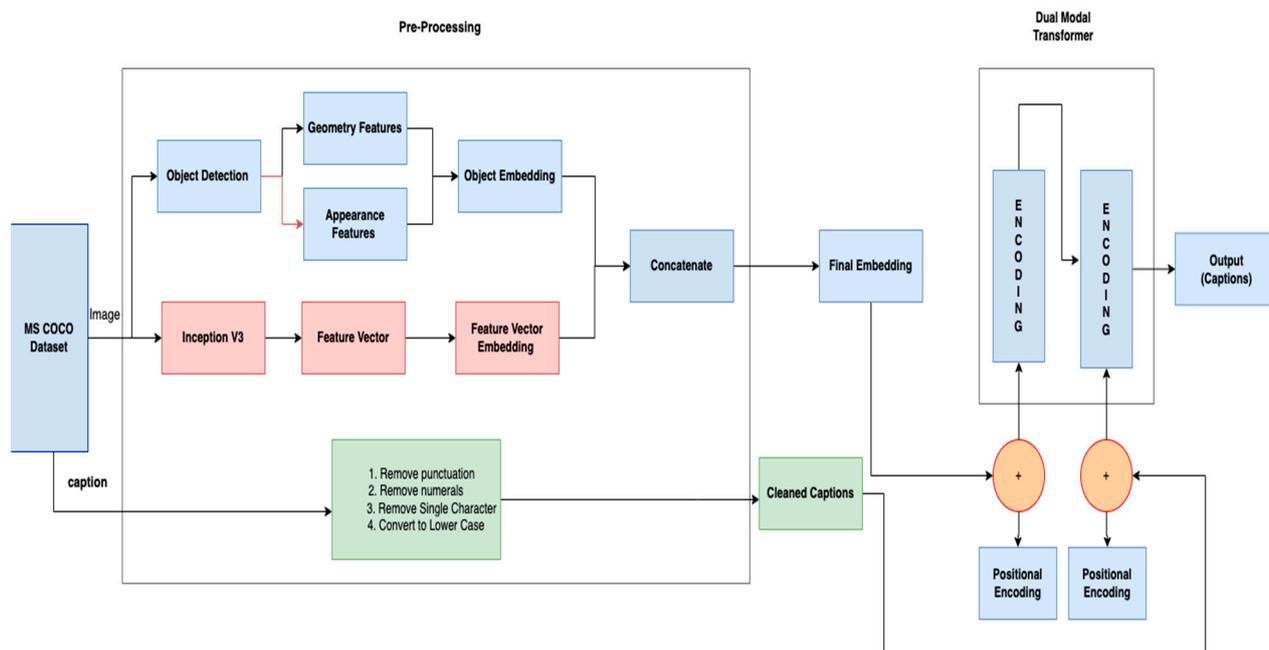
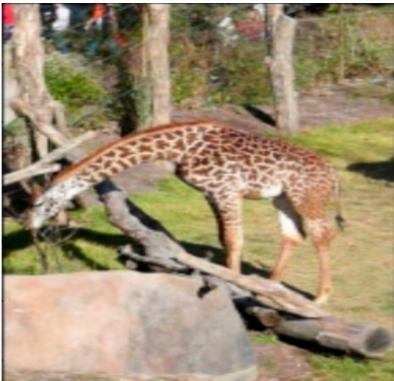


Figure 1. Workflow for Image Captioning Using Dual-Modal Transformer.

3.1. Dataset Description

The dataset used for image captioning is the MS COCO dataset, which contains 330k images and five captions corresponding to each image along with 1.5 object instances [36]. This dataset has been considered one of the most promising datasets in the domain of image captioning as it contains non-iconic images, which makes it different and stands apart from other datasets. The term non-iconic here signifies multiple objects overlapping in the image whereas iconic images are those which constitute a single object. This advantage of the MS COCO dataset becomes very useful while incorporating the labelling task for the images. Table 1 shows some of the images that have been taken from the MS COCO dataset for image captioning.

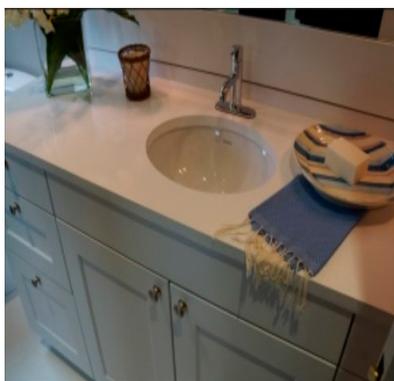
Table 1. Image captioning as given for various images in the Microsoft Common Objects in Context dataset.



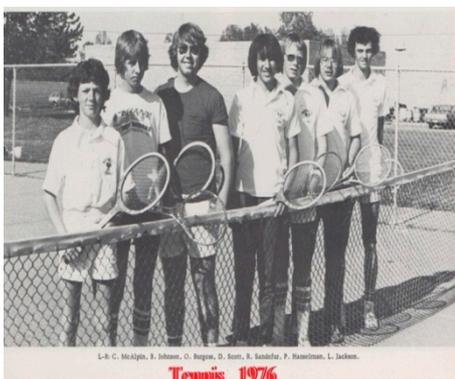
A giraffe standing next to a log in a fenced-in enclosure with people walking on the other side of the fence.
 A giraffe stretching its head over a log to eat weeds.
 A giraffe eating grasses in a wooded area.
 A giraffe in its zoo enclosure bending over to eat.
 A giraffe grazing in the grass next to rocks and trees.



A couple of young men riding bikes down a street.
 A kid wearing a backpack rides his bike near a beach.
 A couple of people are riding bicycles on the beach.
 A young man riding his bike through the park.
 Two kids on bikes with one having a backpack.



Porcelain sink with a white decorated bathroom and blue accents.
 A bathroom sink with flowers and soap on the counter.
 A clean bathroom is seen in this image.
 A clean bathroom vanity artistically displays flowers, soap and a hand towel.
 There is a white sink in the bathroom.



Seven guys holding tennis rackets and standing on a tennis court.
 Seven guys are standing next to the net with their tennis rackets.
 Some guys line up at the net with their rackets.
 A picture of some tennis players playing tennis.
 Some guys holding the tennis racket.

3.2. Pre-Processing

The proposed work primarily aims to create a data frame corresponding to each image consisting of its image id and the respective captions, and generate indexing for each caption related to the image id. Before building a vocabulary of the words in the captions some pre-processing needs to be done before feeding them into the model. This pre-processing is explained in detail in the following sub-sections.

3.2.1. Feature Vector Based on Convolutional Neural Network (CNN)

To extract the feature vectors for the image, the Inception-V3 model has been used excluding the softmax layer. Before feeding the images into the model, the images have been resized to an identical size of 299×299 , which brings the output to the size of $8 \times 8 \times 2048$ for this layer specifically. The fully connected layer at the top has not been included and pre-processing of each image has been initiated with Inception V3. The image features are reshaped to 64×2048 . Tokenisation of the captions has also been incorporated to index-to-word and word-to-index mapping. As all of the captions are of different sizes, padding has been done by feeding zeros to bring all of them to the same length of the longest caption. The vocabulary size has been limited to the top 8000 words. The Inception V3 model is a very popular model for feature extraction in images. The BLEU scores obtained for Inception V3 are the highest as compared to the scores obtained from other models. Furthermore, when models such as Inception V4 are used there has been no improvement in the BLEU scores; however, the time taken has been increased due to more inception modules. Thereby, Inception V3 has been an optimum choice for the feature extraction.

3.2.2. Object Detection

The proposed work models the use of a pre-trained detector called Faster Regions with CNN (RCNN) to detect objects in the image. Faster R-CNN considers the source image, feeding it as an input to a pre-trained CNN model Res-Net, and getting a feature map as the output. This feature map is then passed as an input to another CNN model called region prediction network (RPN) for region proposition.

In a feature map, F-RCNN [37] considers a fixed number of bounding boxes with varied sizes for every spatial position and then incorporates a deep learning method to predict which bounding boxes are most likely to be detected as objects. The detected bounding box for some of the images from the dataset is shown in Figure 2. Post extracting the feature vectors, a region-based convolutional neural network (R-CNN) has been used to extract features in order to form the final image-based embeddings. R-CNN flattens the feature map corresponding to a bounding box followed by a series of fully-connected layers. Two things are achieved as output. The primary one is the score for $n+1$ classes, where n signifies the total number of object classes and $+1$ is the background class for eradicating the bad proposal. The second output has size $4n$ relative to 4-dimensional offset for each of the n classes for decent adjustment of bounding boxes as per the predicted class. The above-mentioned process yields two vectors, one containing the appearance of the detected objects and the other having the location for them. As the final step for object embedding, both the location vectors and appearance vectors have been combined according to Equation (1). This object embedding then served as one of the inputs for the embedding to be fed into the encoder.

$$x_p^{obj} = Lin(W_1 x_p^{bbox}) + Lin(W_2 x_p^{fr}) \quad (1)$$

where, x_p^{obj} is the final object embedding, x_p^{bbox} is the location vector, x_p^{fr} is the appearance vector, W_1, W_2 are linear projected metrics for p -th token and $Lin(\cdot)$ is the layer normalisation function.

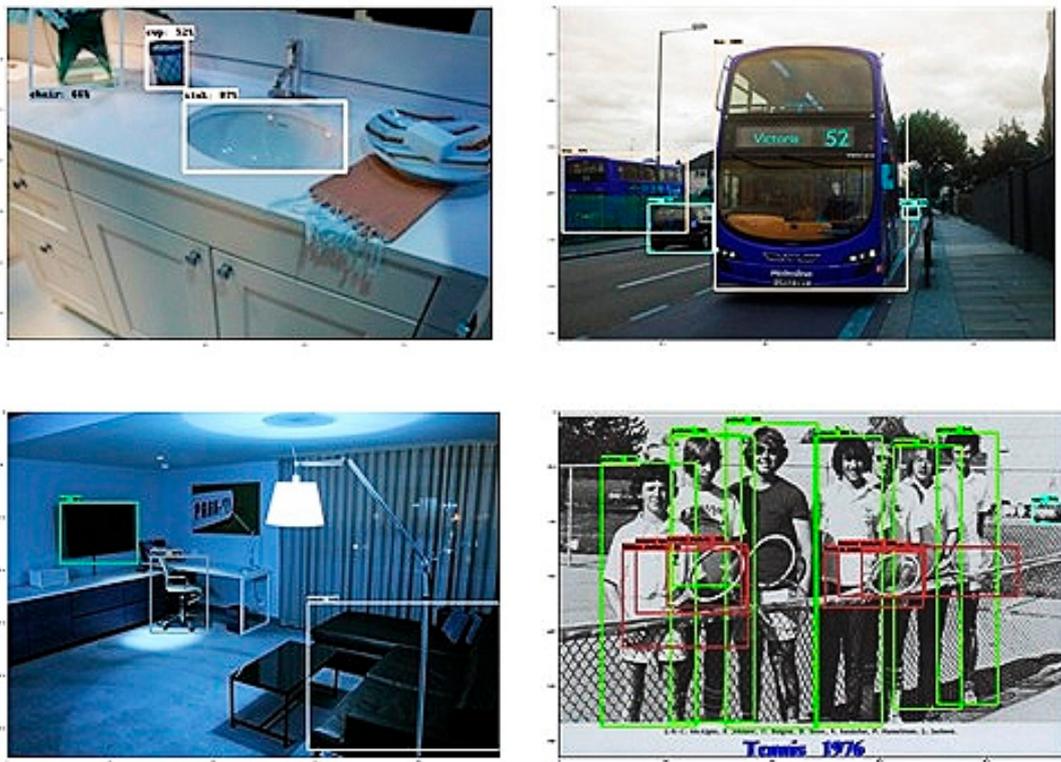


Figure 2. Object detection by F-RCNN and RPN models shown with the bounding boxes.

3.2.3. Data Cleaning

The text cleaning done includes the removal of numeric values, punctuation marks, and single characters. The text is also converted to lowercase before developing the final vocabulary. To mark the beginning and end of a caption, a <start> and <end> token has been added to each one of them, making the model understand the starting and ending points well. The image paths created are the total number of captions that have been incorporated for the image captioning task.

3.3. Positional Encoding

The transformer model promotes the multi-head self-attention mechanism excluding the RNN's recurrence method for accounting the order of the words. The order of words in a sentence has a pivotal role while generating the captions. In the proposed work, each word traverses through the encoder–decoder blocks of the transformer simultaneously; as a result, the model on its own does not take into consideration the position of the words nor does it have the knowledge of each word of the sentence traversing through its stack of blocks. To provide some sense of knowledge to the model about the position of the words in the sentence, positional encoding has been incorporated [38,39]. The encoding taken into consideration creates a vector using the cosine function for every odd index of the input vector, whereas for every even index it creates the vector using the sine function, where both sine and cosine functions are of varied frequencies and depicted in Equations (2) and (3), respectively:

$$PEV_{(posn,2i)} = \sin\left(\frac{posn}{10000^{\frac{2i}{dm}}}\right) \quad (2)$$

$$PEV_{(posn,2i+1)} = \cos\left(\frac{posn}{10000^{\frac{2i}{dm}}}\right) \quad (3)$$

where, PEV is positional encoding vector, i is the index, $posn$ is word's position in the sequence and dm is the size of the embedding vector. This encoding satisfies the criteria that the output is a unique encoding corresponding to the time stamp and the model is capable

of generating long captions with minimal efforts with the consistent distance between any two-time stamps for varied length sentences.

3.4. Dual-Modal Transformer

Transformer networks are very popular in image captioning as these are quite powerful in terms of computational resources during the training and inference; as a result of which, usage at scale is limited for sequences with long-term dependencies [40]. In order to make the transformer networks efficient in the case of very long sentences, an adaptive attention span has been considered where, without drastically increasing the computation time, the attention span has been increased to about 8,000 tokens. The transformer is made up of multiple encoding and decoding units, as shown in its architecture in Figure 3.

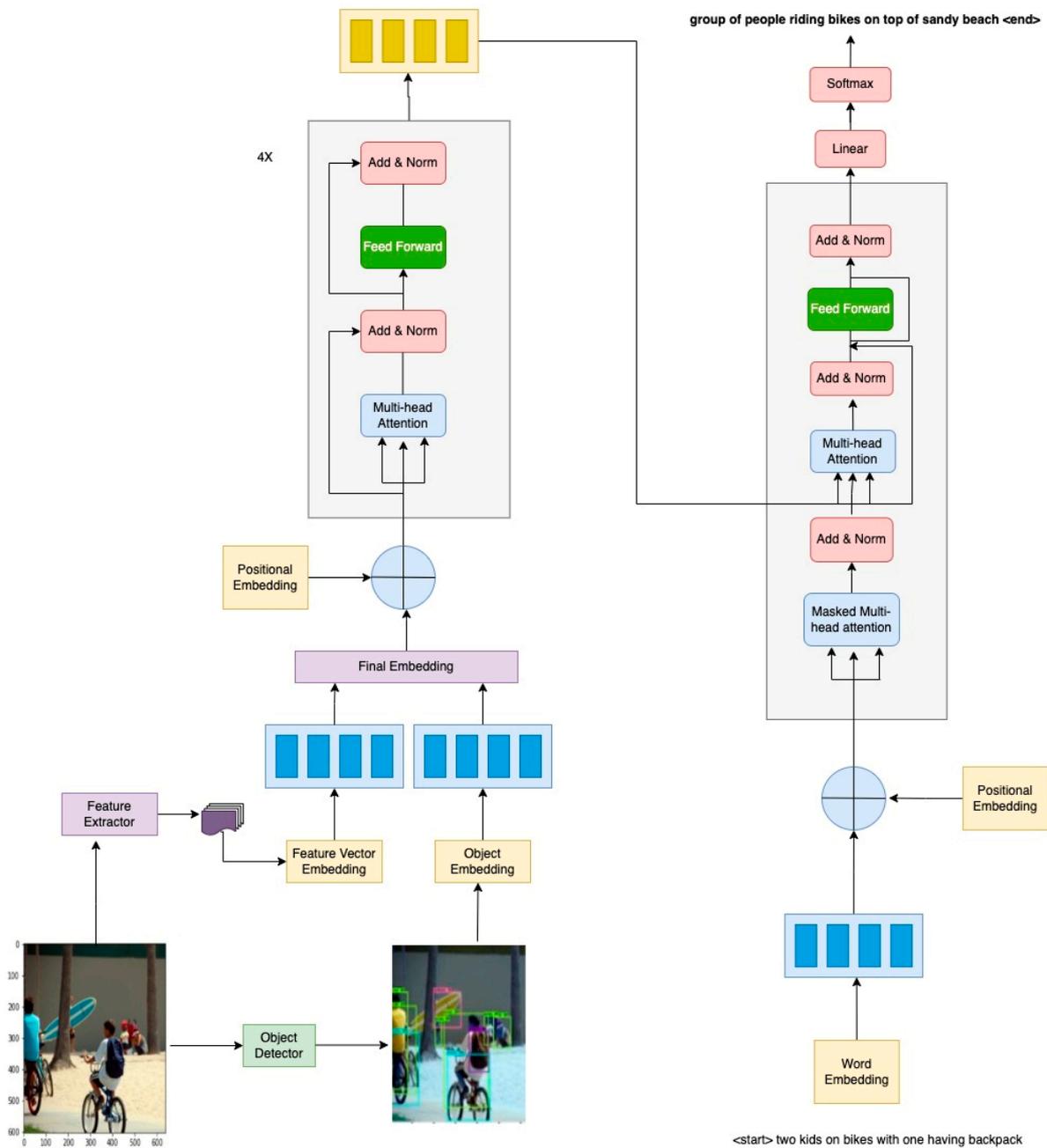


Figure 3. Architecture of the transformer used for image captioning.

3.4.1. Multi-head Attention Layer

The multi-head attention layer in itself is a blend of multiple attention layers. The latter focuses on calculating self-attention using vectors in the form of matrix multiplication. For this, initially, three matrices, namely Query, Key, and Value, are produced for each word by multiplying a previously trained, weight matrix with the input embedding. All of the matrices have their trained weight matrix. Now, the attention function takes these three as inputs and returns attention weights as given in Equation (4).

$$Attention(V, K, Q) = softmax\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (4)$$

For such attention layers, we get attention weight matrices. However, the output of the multi-head attention block needs to be a single matrix. Therefore, all of the attention matrices are integrated and then multiplied with a jointly trained weight matrix. Output of this operation is the final output of the multi-head attention layer. This matrix is then summed with the residual embedding followed by layer normalisation before passing it as an input to the next layer, i.e., feed-forward layer.

3.4.2. Encoder

The input image to the transformer is processed twice: first, to get a set of feature vectors, and then again for a set of detecting objects in the image, which are a crucial part of the caption prediction process, as explained in Section 3.2. Both of these sets are then converted into feature vector embedding and object embedding, respectively. Finally, both the embeddings are combined to form a final embedding vector. This final embedding vector combined with the positional encoding vector acts as the input to the encoder unit. The encoding unit consists of two layers, namely, the multi-head attention layer and the feed-forward layer. A residual association surrounds each of these sub-layers, which is accompanied by layer normalization. The vanishing gradient problem of deep networks can be avoided with these residual connections. In the transformer, such multiple identical units are stacked together, which makes up the complete encoder block. The first unit receives the combined embedding of feature vectors and positional embedding while the rest of the units take the output of their respective previous unit as input. The number of units to be taken is decided experimentally.

3.4.3. Decoder

Like an encoder unit, each decoder unit also consists of sub-layers, namely the masked multi-head attention layer, multi-head attention layer, and feed-forward layer. Both the attention layers work on the same principles at the encoder side. A residual association surrounds each of these sub-layers, followed by layer normalization. The number of such decoder units taking part is equal to the number of encoding units. However, the input of the decoder is variable for the training and testing phases. During training, ground truth caption is fed while predicted words go through the decoder unit in testing. The output obtained after each step in the decoder is part of the input of the next round along with the positional embedding, which indicates the position of a word. The masked multi-head attention layer contains a look ahead as well as a pad mask. As mentioned before, multiple attention units are stacked in the encoder. The output of the topmost unit is converted into a set of attention vectors, which are referred as Key(K) and Value(V). These vectors are fed to the multi-head attention layer of the decoder unit, helping the decoder focus on relevant areas in the input sequence. The output of the stacked decoder is a vector containing float values, which is not the expected result.

For the model to be able to generate a caption for an image, the output should be in the form of words. This is achieved by using a linear layer and a softmax layer. The linear layer converts the output of the decoder into a dimensionally larger vector whose size is equal to the number of unique words in our vocabulary. Now each value in this

newly formed vector corresponds to a unique word in the dictionary. When operated with a softmax layer, as per Equation (5), a score is calculated for each of these values and the word corresponding to the highest score is chosen as the result of that step. The embedding of this word becomes part of the input in the next step and the cycle continues until the end token is obtained.

$$\sigma(\vec{k})_i = \frac{e^{k_i}}{\sum_{j=1}^N e^{k_{ji}}} \tag{5}$$

where,

$\sigma = softmax$

$\vec{k} = inputvector$

$e^{k_i} = standardexponentialfunctionfortheinputvector$

$N = numberofclassesinmulti - classclassifier$

$e^{k_{ji}} = standardexponentialfunctionforoutputvector$

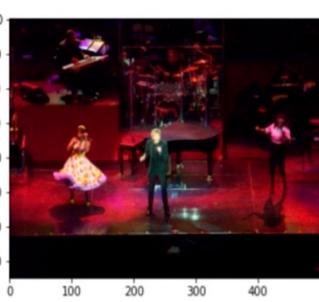
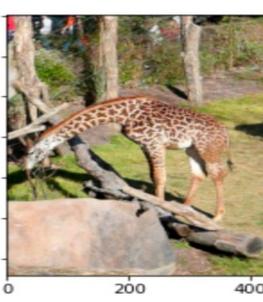
4. Results and Analysis

The transformer architecture has been modelled for generating efficient image captions, taking into account a vocabulary size of 8000 on the MS COCO dataset incorporating 330k images equipped with 1.5 object instances. The efficiency of the model has been compared with two different datasets Flickr 8k, Flickr 30k and the state-of-the-art methods using different captioning metrics. The Flickr 8k and Flickr 30k datasets, as signified by the names, consist of 8000 and 82,783 images, respectively, with five different captions for each image, describing the salient entities and features [16]. Two search methods, beam search by taking the beamwidth of '5', beamwidth '3' and greedy search have been computed on the output probabilities from the model to evaluate the BLEU scores for each dataset, respectively. The captions generated from the transformer model have been listed in Table 2 with the respective images.

Table 2. A comparison of captions generated from greedy search and beam search on Flickr 8k, Flickr 30k and Microsoft Common Objects in Context (MS COCO) datasets.

Flickr 8k	Flickr 30k	MS COCO
		
<p>Real Caption: woman in blue tank top and shorts is jogging down an asphalt road.</p>	<p>Real Caption: football game in progress</p>	<p>Real Caption: two kids on bikes with one having backpack</p>
<p>Greedy Caption: woman jogging up hill</p>	<p>Greedy Caption: two men in white shirts try to tackle player of football players against the team</p>	<p>Greedy Caption: couple of people riding bikes on beach</p>

Table 2. Cont.

Flickr 8k	Flickr 30k	MS COCO
Beam Search (k = 3) Caption: woman in purple shirt and blue shorts is jogging on the path near a mountain	Beam Search (k = 3) Caption: two football players are wearing red and white jerseys.	Beam Search (k = 3) Caption: a group of people riding bikes on top of beach
Beam Search (k = 5) Caption: the woman is wearing purple shirt and blue shorts and jogging on the road in the middle of the dry day	Beam Search (k = 5) Caption: two football players talk during a football game	Beam Search (k = 5) Caption: group of people riding bikes on a top of sandy beach
		
Real Caption: snowboarder flies through the air	Real Caption: here is rod stewart in concert singing live with pat benatar	Real Caption: giraffe grazing in the grass next to rocks and trees
Greedy Caption: snowboarder is jumping in the air after a jump	Greedy Caption: two people are playing the guitars and other people are singing on stage	Greedy Caption: giraffe standing on a dirt road next to trees
Beam Search (k = 3) Caption: the snowboarder is jumping through the air after jumping of the snow-covered mountain jump	Beam Search (k = 3) Caption: a group of people are playing music on stage	Beam Search (k = 3) Caption: there are two giraffes that are standing in the grass
Beam Search (k = 5) Caption: the snowboarder is jumping through the air on their snowboard in front of the ski jump	Beam Search (k = 5) Caption: a group of people are playing musical instruments in what appears to be church	Beam Search (k = 5) Caption: close up view of very pretty giraffes in wooded area

4.1. Greedy Search

Greedy search is a method used in sentence generation in which the word with maximum probability is chosen, then it is set as input for the next step, and this process repeated until the caption length is reached [41]. The proposed model generates the probability of each word at each time step. The algorithm starts with words having maximum probability and then the next word prediction is done greedily, which is used for further computation. Finally, a caption with a relevant set of words has been generated, which has a relatively higher outcome probability [42].

4.2. Beam Search

As observed in the greedy search algorithm, the most likely words are chosen for the output sequence. Contrary to this, another method used is the beam search algorithm, which takes into account all of the possibilities for the sequence and progressively keeps the best ones with it. Now, the major factor in the whole process is the value taken for beam width, say k, which specifies the number of best possibilities to be sent in the next step. In other words, the beam size saves the number of partial hypotheses equal to the unfinished translations beam size in the memory [43]. The number of possibilities explored equals the product of beam width and the vocabulary size.

In the proposed work, beam search operates on a set of 8000 words in the vocabulary with a beam-width of five to predict a caption for the image. Initially, out of 8000 words, the top five most likely first words for the caption are chosen. Then, each of these seven words

is hardwired into the first decoder unit one by one. All of the words from the vocabulary are then evaluated against the first word to find the most likely second word of the caption. Therefore, for each step, a total of 40,000 probabilities are calculated out of them top five are picked and fed into the next iteration. Then, each of these seven pairs is evaluated against all of the vocabulary words to get the third most likely word of the caption. This process goes on until either the end of the caption token is the outcome of the decoder unit or the maximum caption length is reached [44]. The termination can happen at different times for different sets. In the final step, whichever set has the highest probability becomes the predicted caption for that image. The same method has been repeated with beam width three and the results have been recorded. The BLEU scores for the generated captions using the transformer model on Flickr 8k, Flickr 30k and MS COCO datasets have been recorded in Table 3. In the image captioning process, there could be more than one correct caption for an image. Hence, for measuring the accuracy of the model proposed, a metric used conventionally, i.e., BLEU score is used [45]. Bilingual Evaluation Understudy, BLEU, works on the principle of weighted precision in which the words are assigned weights as per their appearance in the reference captions. The same is calculated according to Equation (6). Therefore, if any of the generated captions is close to the reference captions, the value for the BLEU score is high. The reference captions are part of the test set provided by humans as an input for calculating this score.

$$B_n = \frac{\sum_{n\text{-grams} \in y} \text{CountClip}(n\text{-gram})}{\sum_{n\text{-grams} \in y} \text{Count}(n\text{-gram})} \quad (6)$$

where, B_n is the modified precision, n is the number of words taken at a time, y is the predicted caption, $\text{Count}()$ returns the maximum times the n -gram appears in a reference, $\text{CountClip}()$ returns the clipped value of the word by its maximum reference count.

However, the BLEU is dependent not just on the precision but also on a metric called brevity penalty (BP). This helps to balance the short-length captions specifically. This can be calculated as per Equation (7). The value for BP is 1 if the output caption is the same length as any of the reference captions

$$BP = \begin{cases} 1 & \text{if } pred > ref \\ e^{(1 - \frac{ref}{pred})} & \text{if } pred \leq ref \end{cases} \quad (7)$$

where, BP is brevity penalty, $pred$ is word count in predicted caption and ref is word count in reference caption.

Now, based on the number of words used for comparison, BLEU scores are of different types. When contrast is drawn using one word at a time, a BLEU-1 score is obtained. However, when a pair of words, called bigram, is considered at a time, the calculated metric is called the BLEU-2 score. Similarly, trigrams and 4-g help in the computation of the BLEU-3 score and BLEU-4 score, respectively. All of these scores can be calculated by using Equation (8).

$$BLEU_k = BP \cdot \exp\left(\sum_{k=1}^N w_k \log B_n\right) \quad (8)$$

where, k is the number of words, $N = 4$ by default, w_k is the weight for precision, which is $1/N$ hence the default value is 0.25.

Image captioning has been performed on three datasets—Flickr 8k, Flickr 30k, and MS COCO, which have 8000, 30,000, and 82,783 images, respectively. Generally, two patterns have been observed in the BLEU score, either it decreases or it increases from BLEU-1 to BLEU-4. It is observed that the unigram bleu score favours short predictions and when the length of caption predicted is long, the value of the unigram bleu score decreases.

Table 3. Performance comparison in terms of bilingual evaluation understudy scores on Flickr 8k, Flickr 30k and MS COCO datasets with greedy search, beam search with beam width three and beam search with beam width five using the dual-modal transformer model.

Dataset	Greedy Caption				Beam Search (k = 3)				Beam Search (k = 5)			
	B-1	B-2	B-3	B-4	B-1	B-2	B-3	B-4	B-1	B-2	B-3	B-4
Flickr 8k	40.76	32.64	29.97	48.13	50.00	33.96	24.82	31.31	70.47	64.23	51.92	39.22
Flickr 30k	55.55	45.64	42.89	31.55	66.66	57.73	49.39	35.49	71.11	66.69	53.12	41.03
MS COCO	44.12	34.96	49.76	54.50	40.30	44.72	61.70	66.87	85.01	78.37	70.34	48.33

In contrast, the BLEU-4 score is observed to be of high value when the length of the predicted caption is long. From the results, it is evident that beam search predictions are much better than greedy search predictions across all three datasets that we have used. As the size of the dataset increases, the generating ability of the model increases. Teacher forcing has been used for decoder training considering the ground truth captions as input to it at each time step instead of the word that was predicted in the previous time step. The purpose of the same is to increase the speed of training time by a significant value. The proposed work also uses beam search, which selects the word having the highest cumulative score from the total words in its sequence for generating better captions. The maximum BLEU score obtained using a beamwidth of five is 85.01 on the MS COCO dataset as shown in Figure 4d–f.

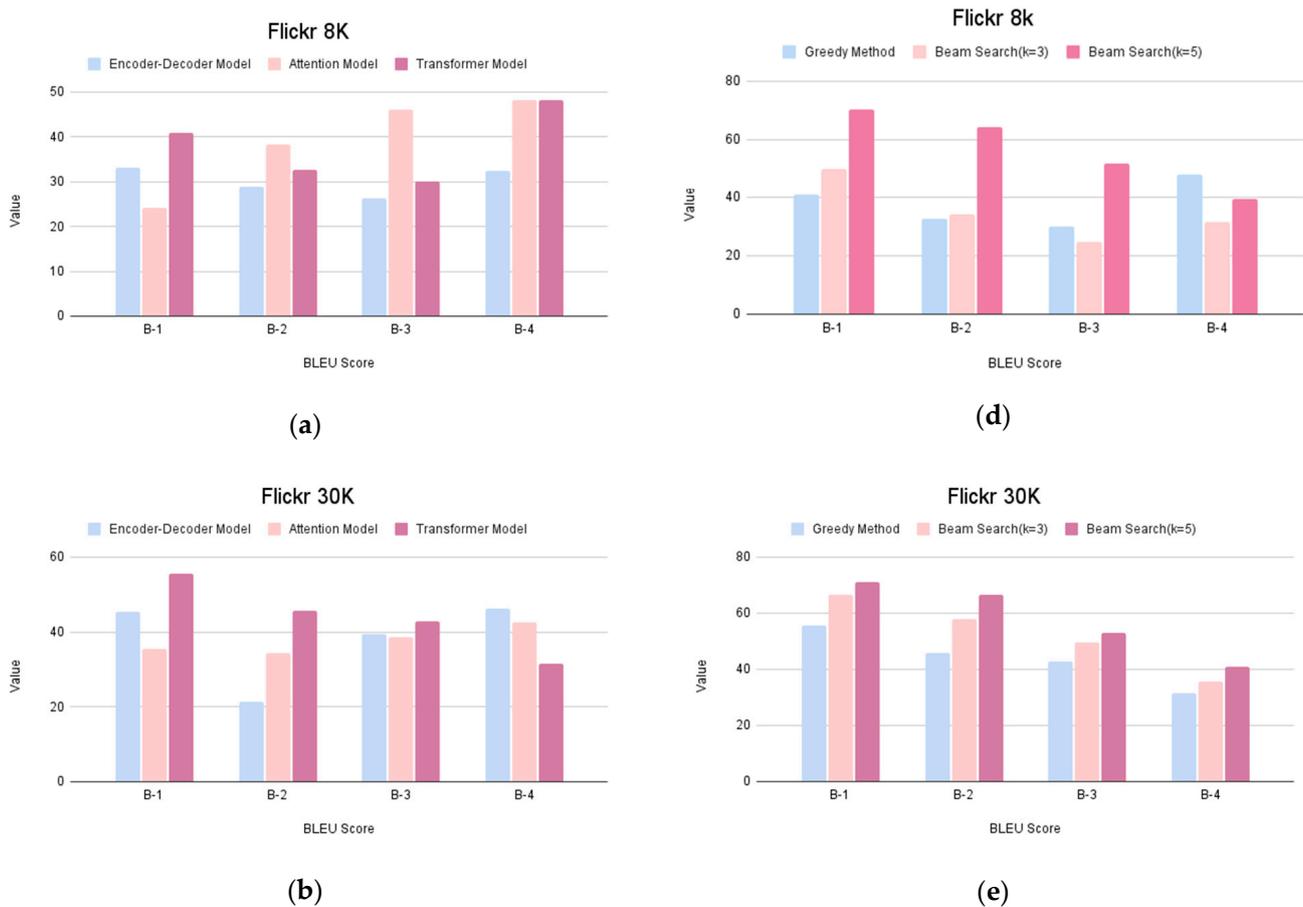


Figure 4. Cont.

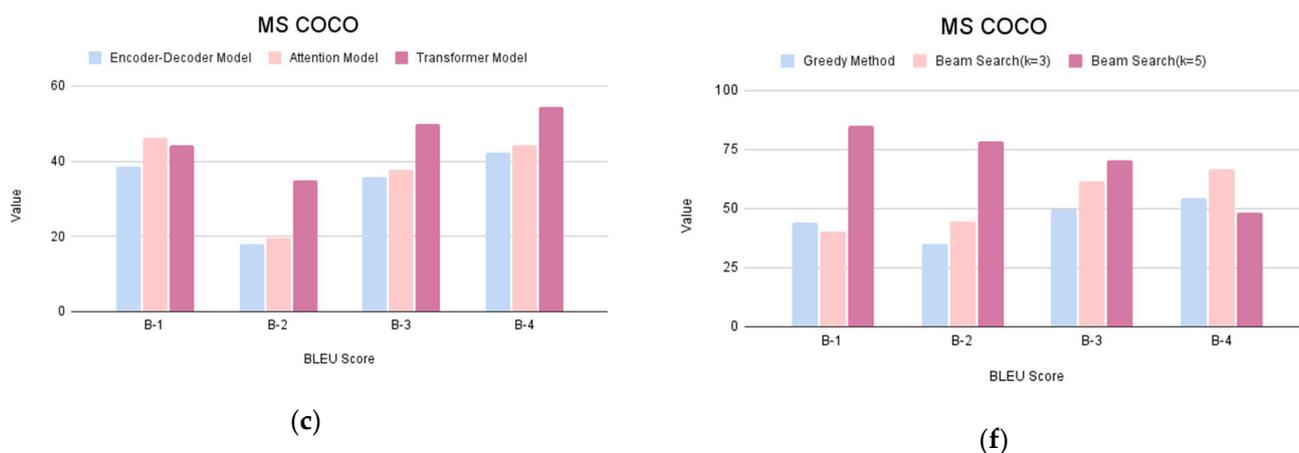


Figure 4. (a–c): shows the BLEU scores of the attention model, encoder–decoder model, and transformer model in which the latter has the highest value of scores; (d–f) shows the BLEU scores of the transformer model using greedy search, beam search ($k = 3$), and beam search ($k = 5$) on Flickr 8k, Flickr 30k, MS COCO datasets.

4.3. Comparative Analysis

Within this implemented work, an attention mechanism in which the given image is divided into n parts primarily, followed by the computation of image representation on each of them, has been incorporated. The generation of a word by RNN has been simultaneously led by the attention mechanism in focussing on relevant portions of a given image, such that the decoder only uses these particular parts of the image. Two models, namely, the attention model and the encoder–decoder model, are applied to three datasets, namely, Flickr8k, Flickr30k, and MS COCO.

4.3.1. Attention Model

The attention model is a divide and conquers processing technique in which the global computations are broken into a set of local computations which is able to process them easily [46]. It develops a context vector using the input sequence. In other words, the input sequence is filtered to make a vector suitable for that particular time step. These context vectors are then processed along with previously predicted words to predict the next target word. In the comparative study, this model is applied to two other datasets apart from the MS COCO dataset, the results of which are shown in Table 4.

4.3.2. Encoder–Decoder Model

The encoder model combines the encoded form of image and text caption and feeds it to the decoder [47]. The proposed model treats CNN as the ‘image model’ and LSTM as the ‘language model’ variable-length text sequences. A merged architecture is then created for training a neural network responsible for handling images and language separately. To encode image features the InceptionV3 model is used and to encode text sequence, a 200-dimensional vector mapping of every word is done on the embedding layer. For caption generation, the greedy search technique is used to choose the best words for captions. The results for the same are listed in Table 5.

BLEU scores have been computed for both models to provide a comparative overview with the transformer model. Table 6 provides the captioning metrics for the encoder–decoder model and attention model on all three datasets, respectively.

Table 4. A comparison of captions generated using the attention model on Flickr 8k, Flickr 30k and Microsoft Common Objects in Context (MS COCO) datasets.

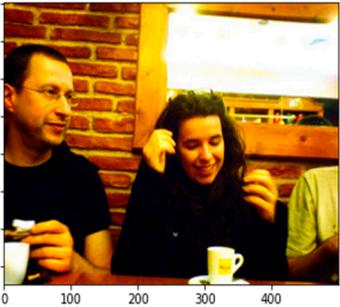
Flickr 8k	Flickr 30k	MS COCO
		
Real Caption: woman in black is sitting next to a man wearing glasses in the coffee shop	Real Caption: woman wearing the striped shirt is looking at mannequin holding up a necklace in a store window	Real Caption: porcelain sink with a white decorated bathroom and blue accents
Predicted Caption: man and woman are sitting down and having to drink	Predicted Caption: woman in a black tank top is holding the cup	Predicted Caption: a couple of a sink and a white towel under a stove

Table 5. A comparison of captions generated using the encoder–decoder model on Flickr 8k, Flickr 30k, and Microsoft Common Objects in Context (MS COCO) datasets.

Flickr 8k	Flickr 30k	MS COCO
		
Real Caption: dog chases ball that has dropped in the water	Real Caption: a young girl in a blue dress is walking on the beach.	Real Caption: a few people riding skissnow-covered snow covered slope
Predicted Caption: two dogs look at something in the water	Predicted Caption: a young boy in a blue shirt is walking on the beach	Predicted Caption: a man standing on top of a snow-covered slope

Table 6. Performance comparison based on different BLEU scores on Flickr 8k, Flickr 30k and MS COCO datasets using the encoder–decoder model and attention model.

Model	Dataset	B-1	B-2	B-3	B-4
Encoder–Decoder Model	Flickr 8k	33.09	28.88	26.37	32.36
	Flickr 30k	45.45	21.32	39.56	46.17
	MS COCO	38.46	17.90	35.62	42.31
Attention Model	Flickr 8k	24.26	38.36	46.07	48.23
	Flickr 30k	35.52	34.24	38.51	42.66
	MS COCO	46.15	19.61	37.62	44.28

From Table 6 it can be observed that the attention model can generate better predictions, and hence this accounts for its better BLEU score than the merged encoder–decoder model. Table 3 shows the BLEU scores of the generated captions evaluated on the greedy search and beam search algorithms with beamwidths of three and five. It is evident from the

captions that the greater the beam width used, the better predictive capability of the model, as observed in Table 2. With a beamwidth of five, the model tries to capture more detail, and hence generate a better caption as compared to the caption generated when using a beamwidth of three.

In case of the transformer model, the input captions can relate to different parts of the image in a better manner, and hence can generate better captions as compared to other models used. The multi-head attention mechanism used in the transformer model enables it to experience different scenes observed in the image at different levels of detail.

4.3.3. Analysis using State-of-the-Art Methods

The performance of the proposed dual-modal transformer is compared with several existing state-of-the-art methods on MS COCO dataset. In the review network [48], RNN-based decoders are used against CNN and RNN-based encoding units for caption prediction while in the adaptive model [49], the model toggles with the choice between attending to the detected objects from the image and choosing visual sentinel. The upcoming CNN architecture combined with the attention mechanism is employed in CNN + Attn. [50]. CPTR [51] exploits sequential raw images for making patches that are fed to both encoder and decoder in different forms. Simulating spatial relationships, the image transformer [35] employs improvised encoder and decoder units. A dual level collaborative transformer (DLCT) approach uses both grid and region features to generate image captions [45]. Researchers introduced the CAAG (context-aware auxiliary guidance) methodology to guide the captioning model, which learns whole semantics by reproducing the latest generation versus global contexts [34]. The encoder–decoder paradigm has been proposed whereas a global-enhanced encoder first encodes the original inputs into highly abstract localised representations and then extracts the intra- and inter-layer global representation. The decoder then implements the suggested global adaptive controller to iteratively incorporate the multimodal information while producing the caption word by word [52]. An Attentive Fourier-Augmented Image Captioning Transformer (AFCT) based methodology has been proposed by the researchers. The main aim of this research is to develop a transformer-based image captioning system which can effectively utilise the information contained in both images and text, with an emphasis on image attributes to generate grammatically, semantically, and syntactically correct captions [53]. The performance of the proposed methodology has been compared and analysed with the state-of-the-art methods. The efficiency and efficacy of the proposed methodology have been compared and analysed using BLUE (B1, B2, B3, B4) score, METEOR [54] and ROUGE [36] values. The result shows that the dual modal transformer achieved the highest BLUE score of 85.01. The detailed comparative analysis has been depicted in Table 7:

Table 7. A comparison of the proposed method with related state-of-the-art techniques.

Technique	B-1	B2	B-3	B-4	METEOR	ROUGE
Review Net [48]	72.0	55.0	41.4	31.3	34.7	68.6
Adaptive [49]	74.8	58.4	44.4	33.6	26.4	55.0
CNN + Attn [50]	71.5	54.5	40.8	30.4	24.6	52.5
CPTR [51]	81.7	66.6	52.2	40.0	29.1	59.2
Multi-modal transformer [5]	81.7	66.8	52.4	40.4	29.4	59.6
Image Transformer [35]	80.8	-	-	39.5	29.1	59.0
DLCT [42]	82.4	67.4	52.8	40.6	29.8	59.8
CAAG [34]	81.1	66.4	51.7	39.6	29.2	59.2
GET [54]	82.1	-	-	40.6	29.8	59.6
AFCT [53]	80.5	-	-	38.7	29.2	58.4
Dual-Modal Transformer (proposed)	85.0	78.4	70.3	48.3	35.4	69.2

4.3.4. Ablation Study

Table 8 shows the ablation study, which lead to the finalisation of the proposed architecture. The study indicates that when the encoder used simple word embeddings as the feature vector, then the BLEU score obtained has been quite low. To improve the BLEU scores, object embeddings based on the features extracted by object detection in the image are combined with word embeddings. This leads to an increase of 7.8 points in the BLEU score. Further, when features extracted by using VGG-16 are combined to features obtained by object embeddings, then the BLEU score is increased by 9.8 points. However, different deep learning models can have different impacts on the BLEU scores based on their architecture and depth. Thereby, VGG-19 has been used along with object and word embeddings, which leads to an increase of 1.3 points in the BLEU scores. On using the Res-Net Model, the BLEU scores further improved by 0.6 points. However, Inception V3 offered the best BLEU score value of 85.0 points. Finally, when Inception V4 has been used, there has been no increase in BLEU score, but complexity has been increased. Thereby, the proposed model has been selected for training the encoder model in the proposed work.

Table 8. Ablation study results used to decide the proposed architecture.

Model	B-1	B2	B-3	B-4	METEOR	ROUGE
Simple Word embeddings	64.5	52.3	51.2	30.3	23.7	52.7
Word Embeddings + Object Embeddings	72.3	63.5	58.6	35.6	27.4	58.3
Word Embeddings + Object Embeddings+ VGG 16 (Feature Embedding)	82.1	72.7	66.7	40.4	30.2	62.6
Word Embeddings + Object Embeddings+ VGG-19 (Feature Embedding)	83.4	76.5	68.3	42.0	32.4	64.5
Word Embeddings + Object Embeddings+ Res-Net (Feature Embedding)	84.0	77.2	69.7	45.4	32.8	66.9
Dual-Modal Transformer with Word Embeddings + Object Embeddings+ InceptionV3 (Feature Embedding) (proposed)	85.0	78.4	70.3	48.3	35.4	69.2
Dual-Modal Transformer with Word Embeddings + Object Embeddings+ InceptionV4 (Feature Embedding)	85.0	78.4	70.3	48.2	35.3	69.1

5. Conclusions and Future Scope

Image captioning tasks pose a bigger challenge for machines as compared to humans. Hence, many methods are proposed to be able to do it efficiently by the former, one of which is the usage of the transformer model as proposed by the authors. The model exploits the advantageous nature of attention block combined with the two methods, i.e., greedy method and beam search, for predicting a caption for a given input image. The proposed work uses three publicly available datasets viz. MS COCO, Flickr 8K and Flickr 30 K. Two different deep networks have been used to extract two different sets of embedding, which are concatenated to form the final embedding. This final embedding is thereby fed into the encoder and decoder to form a final caption.

For measuring the accuracy of the outputs, evaluation metrics, such as BLEU scores, METEOR and ROUGE are used. The BLEU score for greedy and beam search (with beam-width values as three and five) have been evaluated on all three datasets. Furthermore, a comparison is drawn with the existing state-of-the-art models, namely the attention model and encoder–decoder model, which establishes superiority of the proposed model. The BLEU score evaluates to 85.01 for the MS COCO dataset.

However, there are a few areas in which there is a further scope of improvement. The methodology is not colour sensitive and not able to detect a wide range of colours accurately. Furthermore, due to the limited size of vocabulary, the generating capability

of the model is also limited. Further, the model is unable to achieve high precision in the detection of the number of objects present in the image. In the future, these areas can be improved upon. Further, the model can be trained with a larger dataset and the hyper-parameters of the deep networks can be tuned to achieve better BLEU scores.

Author Contributions: Conceptualisation, methodology, and validation, D.K. and V.S.; formal analysis, D.K.; resources, J.D.H. and D.E.P.; data curation, D.K. and V.S.; review and editing, D.K. and V.S.; visualisation, D.K.; project administration, J.D.H. and D.E.P.; funding acquisition, J.D.H. and D.E.P. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Conflicts of Interest: All authors declare that they have no conflicts of interest.

References

1. Hossain, M.Z.; Sohel, F.; Shiratuddin, M.F.; Laga, H. A comprehensive survey of deep learning for image captioning. *ACM Comput. Surv. (CSUR)* **2019**, *51*, 1–36. [\[CrossRef\]](#)
2. Chen, Y.C.; Li, L.; Yu, L.; El Kholy, A.; Ahmed, F.; Gan, Z.; Cheng, Y.; Liu, J. Uniter: Universal image-text representation learning. In Proceedings of the European Conference on Computer Vision, Glasgow, UK, 23–28 August 2020; Springer: Cham, Switzerland, 2020; pp. 104–120.
3. Yao, T.; Pan, Y.; Li, Y.; Mei, T. Exploring visual relationship for image captioning. In Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, 8–14 September 2018; pp. 684–699.
4. Yang, J.; Sun, Y.; Liang, J.; Ren, B.; Lai, S.H. Image captioning by incorporating affective concepts learned from both visual and textual components. *Neurocomputing* **2019**, *328*, 56–68. [\[CrossRef\]](#)
5. Yu, J.; Li, J.; Yu, Z.; Huang, Q. Multimodal transformer with multi-view visual representation for image captioning. *IEEE Trans. Circuits Syst. Video Technol.* **2019**, *30*, 4467–4480. [\[CrossRef\]](#)
6. You, Q.; Jin, H.; Wang, Z.; Fang, C.; Luo, J. Image captioning with semantic attention. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 27–30 June 2016; pp. 4651–4659.
7. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Polosukhin, I. Attention is all you need. *Adv. Neural Inf. Process. Syst.* **2017**, *30*, 1–11.
8. Xu, K.; Ba, J.; Kiros, R.; Cho, K.; Courville, A.; Salakhudinov, R.; Bengio, Y. Show, attend and tell: Neural image caption generation with visual attention. In Proceedings of the International Conference on Machine Learning, Lille, France, 7–9 July 2015; pp. 2048–2057.
9. Karpathy, A.; Joulin, A.; Fei-Fei, L.F. Deep fragment embeddings for bidirectional image sentence mapping. *Adv. Neural Inf. Processing Syst.* **2014**, *27*, 1–9.
10. Li, G.; Zhu, L.; Liu, P.; Yang, Y. Entangled transformer for image captioning. In Proceedings of the IEEE/CVF International Conference on Computer Vision, Seoul, Korea, 27 October–2 November 2019; pp. 8928–8937.
11. Vinyals, O.; Fortunato, M.; Jaitly, N. Pointer networks. In Proceedings of the Advances in Neural Information Processing Systems, Montreal, QC, Canada, 7–12 December 2015; pp. 2692–2700.
12. Wang, J.; Xu, W.; Wang, Q.; Chan, A.B. Compare and reweight: Distinctive image captioning using similar images sets. In Proceedings of the European Conference on Computer Vision, Glasgow, UK, 23–28 August 2020; Springer: Cham, Switzerland, 2020; pp. 370–386.
13. Yeh, R.; Xiong, J.; Hwu, W.M.; Do, M.; Schwing, A. Interpretable and globally optimal prediction for textual grounding using image concepts. *Adv. Neural Inf. Process. Syst.* **2017**, *30*, 1–11.
14. Amirian, S.; Rasheed, K.; Taha, T.R.; Arabnia, H.R. A short review on image caption generation with deep learning. In Proceedings of the International Conference on Image Processing, Computer Vision, and Pattern Recognition (ICIP), Las Vegas, NV, USA, 29 July–1 August 2019; pp. 10–18.
15. Chohan, M.; Khan, A.; Mahar, M.S.; Hassan, S.; Ghafoor, A.; Khan, M. Image Captioning using Deep Learning: A Systematic. *Image* **2020**, 278–286.
16. Nithya, K.C.; Kumar, V.V. A Review on Automatic Image Captioning Techniques. In Proceedings of the 2020 International Conference on Communication and Signal Processing (ICCSPP), Chennai, India, 28–30 July 2020; pp. 0432–0437.
17. Aggarwal, A.; Chauhan, A.; Kumar, D.; Mittal, M.; Roy, S.; Kim, T.H. Video caption based searching using end-to-end dense captioning and sentence embeddings. *Symmetry* **2020**, *12*, 992. [\[CrossRef\]](#)
18. Al-Muzaini, H.A.; Al-Yahya, T.N.; Benhidour, H. Automatic arabic image captioning using RNN-LSTM-based language model and CNN. *Int. J. Adv. Comput. Sci. Appl.* **2018**, *9*, 67–73. [\[CrossRef\]](#)
19. Herdade, S.; Kappeler, A.; Boakye, K.; Soares, J. Image captioning: Transforming objects into words. *Adv. Neural Inf. Process. Syst.* **2019**, 1–11.
20. Cornia, M.; Stefanini, M.; Baraldi, L.; Cucchiara, R. Meshed-memory transformer for image captioning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Seattle, WA, USA, 13–19 June 2020; pp. 10578–10587.

21. Fei, Z. Attention-Aligned Transformer for Image Captioning. In Proceedings of the Association for the Advancements of Artificial Intelligence, Virtual, 22 February–1 March 2022.
22. Yan, S.; Xie, Y.; Wu, F.; Smith, J.S.; Lu, W.; Zhang, B. Image captioning via hierarchical attention mechanism and policy gradient optimization. *Signal Process.* **2020**, *167*, 107329. [[CrossRef](#)]
23. Chen, Z.; Huang, S.; Tao, D. Context refinement for object detection. In Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, 8–14 September 2018; pp. 71–86.
24. Ng, E.G.; Pang, B.; Sharma, P.; Soricut, R. Understanding Guided Image Captioning Performance across Domains. *arXiv* **2020**, arXiv:2012.02339.
25. Chen, L.; Zhang, H.; Xiao, J.; Nie, L.; Shao, J.; Liu, W.; Chua, T.S. SCA-CNN: Spatial and channel-wise attention in convolutional networks for image captioning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 5659–5667.
26. Atliha, V.; Šešok, D. Comparison of VGG and ResNet used as Encoders for Image Captioning. In Proceedings of the 2020 IEEE Open Conference of Electrical, Electronic and Information Sciences (ESTREAM), Vilnius, Lithuania, 30 April 2020; pp. 1–4.
27. Zhou, L.; Xu, C.; Koch, P.; Corso, J.J. Watch what you just said: Image captioning with text-conditional attention. In Proceedings of the Thematic Workshops of ACM Multimedia, Mountain View, CA, USA, 23–27 October 2017; pp. 305–313.
28. Iwamura, K.; Louhi Kasahara, J.Y.; Moro, A.; Yamashita, A.; Asama, H. Image Captioning Using Motion-CNN with Object Detection. *Sensors* **2021**, *21*, 1270. [[CrossRef](#)] [[PubMed](#)]
29. Aneja, J.; Deshpande, A.; Schwing, A.G. Convolutional image captioning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–23 June 2018; pp. 5561–5570.
30. Lu, J.; Xiong, C.; Parikh, D.; Socher, R. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 375–383.
31. Ji, J.; Luo, Y.; Sun, X.; Chen, F.; Luo, G.; Wu, Y.; Gao, Y.; Ji, R. Improving image captioning by leveraging intra- and inter-layer global representation in transformer network. In Proceedings of the AAAI Conference on Artificial Intelligence, Virtual, 2–9 February 2021; Volume 35, pp. 1655–1663.
32. Deshpande, A.; Aneja, J.; Wang, L.; Schwing, A.G.; Forsyth, D. Fast, diverse and accurate image captioning guided by part-of-speech. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Long Beach, CA, USA, 15–20 June 2019; pp. 10695–10704.
33. Anderson, P.; Fernando, B.; Johnson, M.; Gould, S. Guided open vocabulary image captioning with constrained beam search. *arXiv* **2016**, arXiv:1612.00576.
34. Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Chintala, S. Pytorch: An imperative style, high-performance deep learning library. *Adv. Neural Inf. Process. Syst.* **2019**, *32*, 1–12.
35. Shukla, N.; Fricklas, K. *Machine Learning with TensorFlow*; Manning: Greenwich, UK, 2018.
36. Lin, T.Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Gao, Y.; Zitnick, C.L. Microsoft coco: Common objects in context. In Proceedings of the European Conference on Computer Vision, Zurich, Switzerland, 6–12 September 2014; Springer: Cham, Switzerland, 2014; pp. 740–755.
37. Girshick, R. Fast R-CNN. In Proceedings of the IEEE International Conference on Computer Vision, Santiago, Chile, 7–13 December 2015; pp. 1440–1448.
38. Yang, Z.; Zhang, Y.J.; Huang, Y. Image captioning with object detection and localization. In Proceedings of the International Conference on Image and Graphics, Shanghai, China, 13–15 September 2017; Springer: Cham, Switzerland, 2017; pp. 109–118.
39. Parameswaran, S.N.; Das, S. A Bottom-Up and Top-Down Approach for Image Captioning using Transformer. In Proceedings of the 11th Indian Conference on Computer Vision, Graphics and Image Processing, Hyderabad, India, 18–22 December 2018; pp. 1–9.
40. Luo, Y.; Ji, J.; Sun, X.; Cao, L.; Wu, Y.; Huang, F.; Gao, Y.; Ji, R. Dual-level collaborative transformer for image captioning. *arXiv* **2021**, arXiv:2101.06462.
41. Jiang, W.; Ma, L.; Jiang, Y.G.; Liu, W.; Zhang, T. Recurrent fusion network for image captioning. In Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, 8–14 September 2018; pp. 499–515.
42. Melas-Kyriazi, L.; Rush, A.M.; Han, G. Training for diversity in image paragraph captioning. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, 31 October–4 November 2018; pp. 757–761.
43. Ranzato, M.A.; Chopra, S.; Auli, M.; Zaremba, W. Sequence level training with recurrent neural networks. *arXiv* **2015**, arXiv:1511.06732.
44. Biswas, R.; Barz, M.; Sonntag, D. Towards explanatory interactive image captioning using top-down and bottom-up features, beam search and re-ranking. *KI-KünstlicheIntelligenz* **2020**, *34*, 571–584. [[CrossRef](#)]
45. Papineni, K.; Roukos, S.; Ward, T.; Zhu, W.J. Bleu: A method for automatic evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics, Philadelphia, PA, USA, 6–12 July 2002; pp. 311–318.
46. Chen, H.; Ding, G.; Lin, Z.; Zhao, S.; Han, J. Show, Observe and Tell: Attribute-driven Attention Model for Image Captioning. In Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, Stockholm, Sweden, 13–19 July 2018; pp. 606–612.
47. Xiao, X.; Wang, L.; Ding, K.; Xiang, S.; Pan, C. Deep hierarchical encoder–decoder network for image captioning. *IEEE Trans. Multimed.* **2019**, *21*, 2942–2956. [[CrossRef](#)]

48. Yang, Z.; Yuan, Y.; Wu, Y.; Cohen, W.W.; Salakhutdinov, R.R. Review networks for caption generation. *Adv. Neural Inf. Processing Syst.* **2016**, 1–9.
49. Liu, W.; Chen, S.; Guo, L.; Zhu, X.; Liu, J. Cptr: Full transformer network for image captioning. *arXiv* **2021**, arXiv:2101.10804.
50. He, S.; Liao, W.; Tavakoli, H.R.; Yang, M.; Rosenhahn, B.; Pugeault, N. Image captioning through image transformer. In Proceedings of the Asian Conference on Computer Vision, Kyoto, Japan, 30 November–4 December 2020.
51. Song, Z.; Zhou, X.; Mao, Z.; Tan, J. Image captioning with context-aware auxiliary guidance. In Proceedings of the AAAI Conference on Artificial Intelligence, Virtual, 2–9 February 2021; Volume 35, pp. 2584–2592.
52. Osolo, R.I.; Yang, Z.; Long, J. An Attentive Fourier-Augmented Image-Captioning Transformer. *Appl. Sci.* **2021**, *11*, 8354. [[CrossRef](#)]
53. Banerjee, S.; Lavie, A. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization, Ann Arbor, MI, USA, 29 June 2005; pp. 65–72.
54. Lin, C.Y. Rouge: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*; Association for Computational Linguistics: Barcelona, Spain, 2004; pp. 74–81. Available online: <https://aclanthology.org/W04-10> (accessed on 27 April 2022).