# Combining NLP and Deep Learning Techniques to Generate Captions

*Leo Francis M[1], Darshan K S[2], Ankith M C[3], Divakara V[4]*

[1,2,3,4]Presidency University Bangalore 560047 India
Leof5791@gmail.com
DOI - https://doi.org/10.55248/gengpi.4.523.42704

ABSTRACT

Picture captioning ranges the areas of computer vision and common dialect handling. The picture captioning errand generalizes protest location where the portrayals are a single word. As of late, most investigate on picture captioning has centred on profound learning procedures, particularly Encoder-Decoder models with Convolutional Neural Organize (CNN) highlight extraction. Be that as it may, few works have attempted utilizing question location highlights to extend the quality of the created captions. This paper presents an attention-based, Encoder-Decoder profound engineering that produces utilize of convolutional highlights extracted from a CNN demonstrate pre-trained on ImageNet (Xception), at the side question highlights extricated from the YOLOv4 show, pre-trained on MS COCO. This paper moreover presents a modern positional encoding plot for protest highlights, the "importance factor." Our demonstrate was tried on the MS COCO and Flickr30k datasets, and the execution is compared to execution in comparative works. Our unused include extraction plot raises the CIDEr score by 15.04%.

KEYWORDS: Image Captioning, Yolo, Object Detection, Image Processing, NLP, Deep Learning, CNN, Attention, Decoder, Encoder.

## INTRODUCTION

The subject of independently creating expressive sentences for pictures has invigorated intrigued in normal dialect handling and computer vision investigate in recent years. Image captioning could be a key errand that requires a semantic comprehension of pictures as well as the capacity to create precise and exact depiction sentences.

Within the huge information time, pictures are one of the foremost accessible information sorts on the Web, and the require for commenting on and labelling them expanded. In this way, picture captioning frameworks are an illustration of huge information issues as they center on the volume viewpoint of huge information. For case, the MS COCO dataset contains around 123,000 pictures (25 GB). This includes the requirement of effectively utilize of assets and cautious plan of tests.

Early approaches utilized format strategies, attempting to fill predefined layouts of content by features extricated from pictures. Current frameworks advantage from the accessible computing control and utilize profound learning procedures.

One of the foremost effective strategies for picture captioning is actualizing an Encoder-Decoder design. It encodes pictures to a high-level representation at that point translates this representation employing a dialect era demonstrate, like Long Short-Term Memory (LSTM), Gated Repetitive Unit (GRU) or one of their variations.

The consideration component has illustrated its viability in sequence-to-sequence applications, particularly picture captioning and machine interpretation. It increases accuracy by constraining the demonstrate to concentrate on the critical parts of the input when producing yield groupings.

To get it an image, many present-day profound learning models utilize existing pre-trained Convolutional Neural Systems (CNNs) to extricate frameworks of highlights from the final convolutional layers. This makes a difference to get a handle on numerous viewpoints of the objects and their connections within the picture and represent the picture at a better level.

As of late, a few works attempted to utilize protest highlights for picture captioning. Among the models utilized are the YOLOv3, YOLOv4 and YOLO9000, which are known for their speed, exactness, and viability for real-time applications. Question highlights are often a cluster of question labels, where each question tag contains the bounding box information, question course and certainty rate. This work explores the theory that exploiting such highlights might increment exactness in picture captioning which utilizing all protest highlights makes a difference to precisely mirror the human visual understanding of scenes. This paper points to display a demonstrate that makes utilize of this sort of highlights through a straightforward engineering and assess the comes about.

Segment two of the paper will handle the related works within the space. Area three presents our strategy, which incorporates the proposed show and the pre-processing that was performed on the information. In segment four, the tests plan and results are expounded, and a comparison to past works is appeared. Area five concludes the paper and presents plans for future works.

## RESEARCH METHODOLOGY

The exploratory strategy includes extricating question highlights from the YOLO demonstrate and presenting them at the side CNN convolutional highlights to a straightforward profound learning demonstrate that employments the far-reaching Encoder-Decoder engineering with the consideration instrument. "Results and discussion" segment compare the distinction in comes about some time recently and after including the protest highlights. Although past investigate encoded protest highlights as a vector, we include question highlights in a straightforward concatenation way and accomplished a great advancement. We too test the effect of sorting the question labels extricated from YOLO concurring to a metric that we propose here.

## DATASETS UTILIZED

We test our strategy on two datasets utilized often for picture captioning:

MS COCO and Flickr30k. Table 1 contains a brief comparison between them. They are both collected from the Flickr photo sharing site and comprise of real-life pictures, clarified by people (five comments per picture).

## EVALUATION METRICS

We utilize a set of assessment measurements that are broadly utilized within the picture captioning field. BLEU measurements are commonly utilized in robotized content assessment and measure the correspondence between a machine interpretation yield and a human interpretation; within the case of picture captioning, the machine interpretation yield compares to the naturally delivered caption, and the human interpretation compares to the human portrayal of the picture. METEOR is computed utilizing the consonant cruel of unigram precision and review, with the review having a better weight than the exactness, as takes after:

$$METEOR=[10*Precision*Recall/(Recall+9*Precision)]$$

ROUGE-L employments a Longest Common Subsequence (LCS) score to evaluate the ampleness and familiarity of the created content, whereas CIDEr centres on grammaticality and saliency. Zest assesses the semantics of the created text by creating a "scene graph" for both the initial and created captions, and after that as it were matches the terms if their lemmatized WordNet representations are indistinguishable. BLEU, METEOR, and ROUGE have moo relationships with human quality tests, whereas Zest and CIDEr have a distant better relationship but are more troublesome to optimize.

## MODEL

Our demonstrate employments an attention-based Encoder-Decoder engineering. It has two strategies of highlight extraction for picture captioning:

a picture classification CNN (Xception), and a protest discovery demonstrate (YOLOv4). The yields of these models are combined by concatenation to deliver a include lattice that carries more data to the dialect decoder to anticipate more exact depictions. Not at all like others' works that implanted question highlights some time recently combining them with CNN highlights, we utilize crude protest format data specifically. Dialect era is done utilizing a consideration module (Bahdanau consideration), a GRU and two completely associated layers. Our demonstrate is basic, quick to prepare and assess, and creates captions utilizing consideration.

We accept that on the off chance that people can advantage from question highlights (such as the course of question, its position, and estimate) to superior get it a picture, a computer demonstrate can advantage from this data as well. A scene containing a bunch of individuals standing near together, for illustration, may propose a assembly, while inadequate swarms can show an open area. Figure 1 delineates our demonstrate.
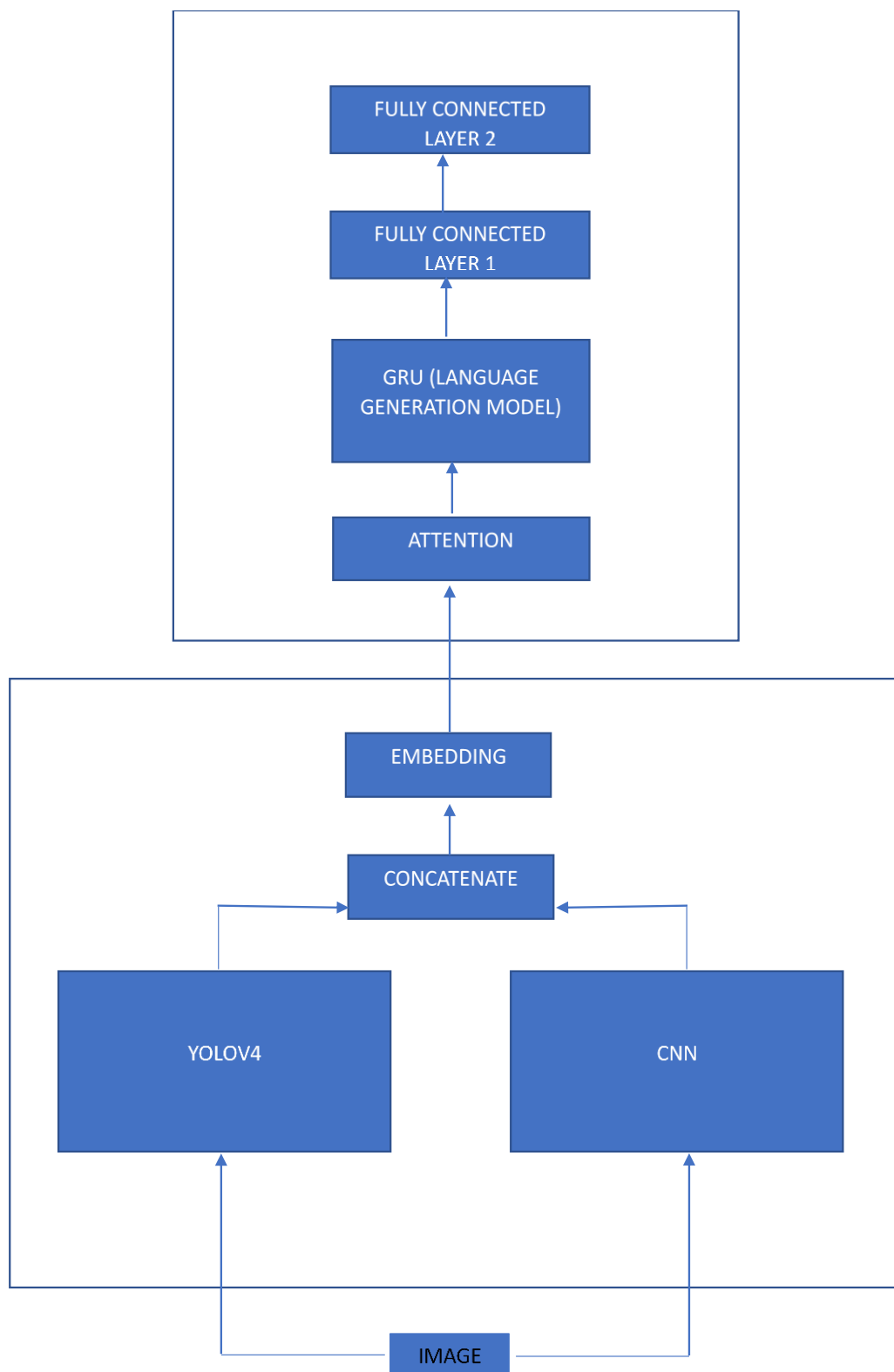
Figure: 1

## IMAGE ENCODING

### A. PRE-TRAINED IMAGE CLASSIFICATION CNN

In this work, we utilize the Xception CNN pre-trained on ImageNet to extricate spatial highlights.

Xception (Extraordinary form of Beginning) is propelled by Beginning V3, but rather than Initiation modules, it has 71 layers with an adjusted depth-wise distinguishable convolution. It outflanks Initiation V3 much obliged to demonstrate parameter utilization way better.

We extricate highlights from the final layer some time recently the completely associated layer, taking after later works in picture captioning. This permits the by and large show to pick up insight about the objects within the picture and the connections between them rather than fair centring on the picture course.

In a past work, distinctive highlight extraction CNN models have been compared for picture captioning applications. The comes about appeared that Xception was among the foremost strong in extricating highlights and for this it was chosen as the include extraction show in this think about. The yield of this arrange is of shape ($10 \times 10 \times 2048$), which was squashed to ($100 \times 2048$) for ease of lattice dealing with.

## B. OBJECT DETECTION MODEL

Our strategy employments the YOLOv4 demonstrate since of its speed and great exactness, which make it reasonable for enormous information and real-time applications. The extricated highlights are a list of protest highlights, with each question include containing the X facilitate, Y facilitate, width, stature, certainty rate (from to 1 comprehensive), course number and a novel discretionary "importance factor".

Taking after human instinct, frontal area objects are ordinarily bigger and more imperative when depicting a picture, and foundation objects are regularly littler and less critical. Moreover, it makes sense to utilize more precise pieces of data than to utilize less precise ones. Thus, our significance figure tries to adjust the significance of the frontal area huge objects and objects with tall certainty rates. The equation to calculate it for a single protest is as takes after:

Significance Factor=Confidence Rate $\times$ Object Width $\times$ Object Height

The significance figure gives a better score to closer view huge objects over foundation little ones, and higher score to objects with a tall certainty over objects with less certainty.

After extricating question highlights, the significance figure is calculated for each protest and concatenated to its tag. At that point, all objects in the list are sorted agreeing to this significance calculate utilizing the speedy sort calculation. Not at all like past works, our strategy makes utilize of all the image's protest data. Since of the measure confinement within the yield of the CNN, we utilize up to 292 objects, each with seven traits (counting the significance calculate), which is ordinarily sufficient to speak to vital objects in a picture.

The list of highlights is straightened into a 1D cluster, of length less than 2048. It is at that point cushioned with zeros to length 2048 to be consistent with the yield of the CNN module. The yield of this organize is a cluster ($1 \times 2048$).

As for calculating the certainty score, YOLO partitions a picture into a lattice. B bounding boxes and certainty scores for these boxes are anticipated in each of these network cells. The certainty score indicates how sure the show is that the box incorporates a protest, as well as how exact the show accepts the box that anticipated is. The question discovery calculation is assessed utilizing Crossing point over Union (IoU) between the anticipated box and the ground truth. It analyses how comparable the anticipated box is to the ground truth by calculating the cover between the ground truth and the anticipated bounding box. A cell's certainty score ought to be zero on the off chance that no question exists in there. The equation for calculating the confidence score is:

C=Pr(object)*IoU

## C. CONCATENATION AND EMBEDDING

In arrange to require advantage of the picture classification highlights and the question discovery highlights, we include this concatenation step, where we connect the yield of the YOLOv4 subsystem as the final push within the yield of arrange 1. The yield of this organize is of shape ($101 \times 2048$).

The implanting is done utilizing one completely associated layer of length 256. This arrange guarantees a steady estimate of the highlights and maps the include space to a littler space fitting for the dialect decoder.

## D. ATTENTION

Our strategy employments the Bahdanau delicate consideration framework. This deterministic consideration component makes the demonstrate as a entire smooth and differentiable.

The term "attention" alludes to a methodology that mimics cognitive consideration. The impact highlights the foremost vital parts of the input information whereas blurring the rest. The concept is that the organize ought to devote more noteworthy computer assets to that little but basic parcel of the information. Which component of the information is more important than others is decided by the setting and is learned by angle plummet utilizing preparing information. Characteristic dialect preparing and computer vision utilize consideration in a few machine learning errands.

The attention mechanism was made to extend the execution of the encoder-decoder engineering for machine interpretation. And as picture captioning can be seen as a particular case of machine interpretation, consideration demonstrated valuable when analysing pictures as well. The consideration instrument was intended to permit the decoder to utilize the foremost pertinent parts of the input grouping in a adaptable way by combining all of the encoded input vectors into a weighted combination, with the foremost significant vectors accepting the most elevated weights.

Consideration takes after the human instinct of centring on diverse parts of a picture when portraying it. Utilizing protest discovery highlights too takes after the instinct that knowing approximately question classes and positions aid get a handle on more almost the picture than simple convolutional features. When consideration is utilized to both highlight sorts, the framework will center on distinctive highlights of both protest classes and positions within the same picture. Figure 2 delineates utilizing attention for picture captioning.
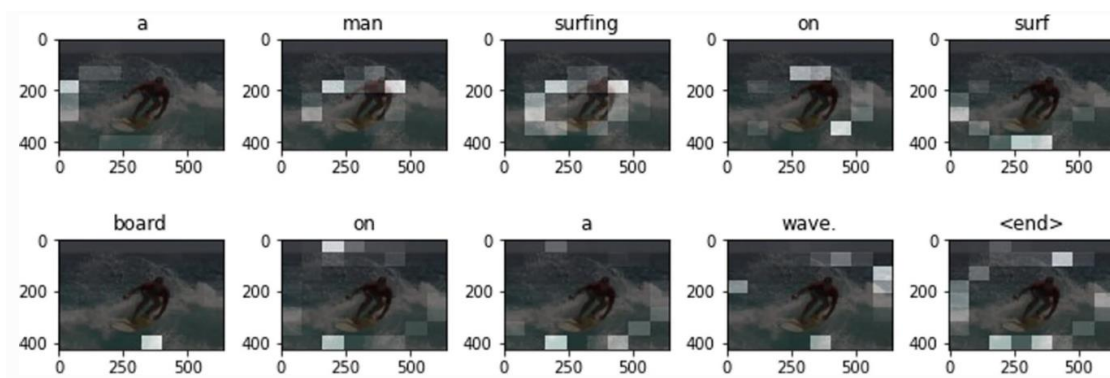


Figure: 2

## LANGUAGE DECODER

For interpreting, a GRU is utilized to abuse its speed and moo memory utilization. It produces a caption by producing one word at each time step, conditioned on a setting vector, the past covered up state, and the already produced words. The show is prepared utilizing the backpropagation calculation deterministically.

The GRU is taken after by two completely associated layers. The primary one is of length 512, and the moment one is of the measure of the lexicon to create yield content.

The preparing prepare for the decoder is as follows:

1. The highlights are extricated at that point passed through the encoder.

2. The decoder gets the encoder yield, covered up state (initialized to 0), and decoder input (which is the begin token).

3. The decoder returns the forecasts as well as the covered-up state of the decoder.

4. The covered-up state of the decoder is at that point passed back into the show, and the misfortune is calculated utilizing the forecasts.

5. To decide another decoder input, "teacher forcing" is employed, which could be a method that passes the target word as another input to the decoder.

## PRE-PROCESSING

This area presents the pre-processing calculation that was performed on the information:

1. Sort the dataset at irregular into image-caption sets. This makes a difference the preparing prepare to meet quick and anticipates any predisposition amid the preparing. Hence, anticipating the demonstrate from learning the arrange of preparing.

2. Perused and interpret the pictures.

3. Resize the pictures to the CNN necessities:  anything the measure of the picture is, it is resized to 299 × 299 as required by the Xception CNN show.

4. Tokenization of the content. Tokenization breaks the crude content into words, that are isolated by accentuations, extraordinary characters, or white spaces. The separators are disposed of.

5. Number the tokens, sort them by recurrence and select the beat 15,000 most common words as the system's lexicon. This dodges over-fitting by disposing of terms that are not likely to be valuable.

6. Create word-to-index and index-to-word structures. They are at that point utilized to interpret token arrangements into word identifier arrangements.

7. Cushioning. As sentences can be diverse in length, we got to have the inputs with the same measure, typically where the cushioning is essential. Here, identifier groupings are cushioned at the conclusion with invalid tokens to guarantee that they are all the of same length.

**RESULTS AND DISCUSSION**

Our code is composed within the Python programming dialect utilizing TensorFlow library The CNN usage and prepared demonstrate were imported from Keras library, and a YOLOv4 demonstrate pre-trained on MS COCO was imported from the yolov4 library. This work employments the MS COCO assessment instrument to calculate scoresFootnote4.

Tests are conducted on two broadly utilized datasets for picture caption era:

MS COCO and Flickr30k. Each picture has five reference captions in these two datasets, which contain 123,000 and 31,000 pictures, individually. For MS COCO, 5000 pictures are saved for approval and 5000 pictures are saved for checking concurring to Karpathy's part. Within the case of Flicker30k dataset, 29,000 pictures are utilized for arrangement, 1000 for approval, and 1000 for testing. The demonstrate was prepared for 20 epochs and utilized Sparse Categorical Cross Entropy as the misfortune work. For the optimizer, Adam optimizer was utilized.

In arrange to subjectively compare the printed yields of the approach, we show in below figures a subjective comparison between the comes about with protest highlights and without them. We take note that the contrast is surprising, and the expansion of the protest highlights makes the sentences more striking syntactically, and with less question botches. In Fig:3 for illustration, a skier was distinguished rather than fair the skiing boots. In Fig:4, the demonstrate some time recently consolidating question highlights had blended up individuals and snow sheets. In Fig:5, the two dairy animals were accurately distinguished after including protest highlights. In Fig:6, The show without question highlights falsely identified a man within the picture. In Fig:7, the demonstrate might not recognize the third bear without question highlights. In Fig:8, question features helped to distinguish a gather of individuals rather than as it were two ladies.



**Baseline model:** this is up in snow pants jumping on a big snowy mountain at night.
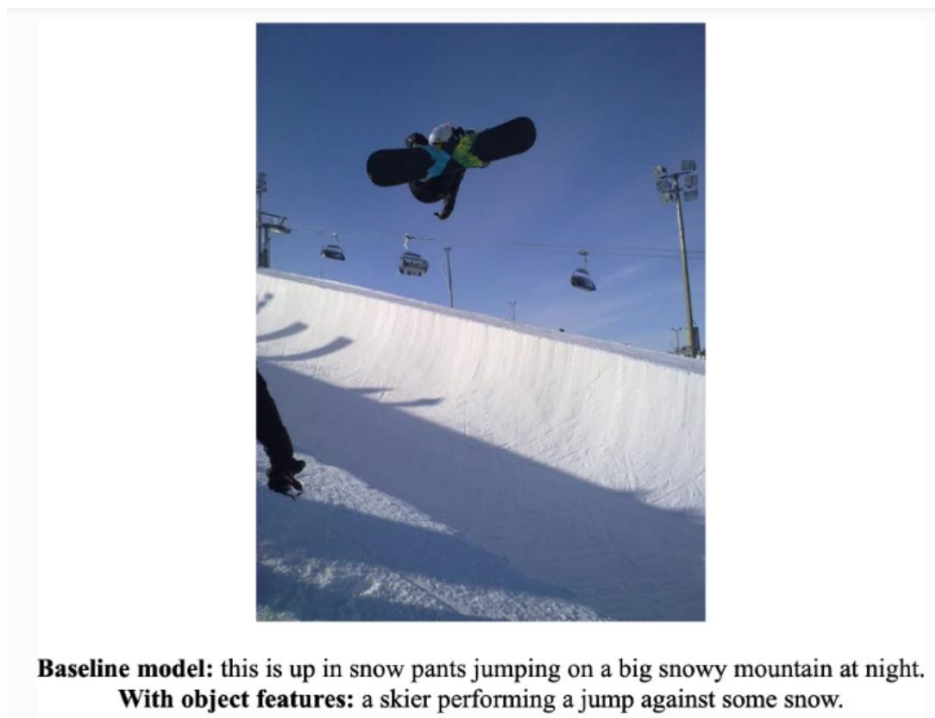**With object features:** a skier performing a jump against some snow.

Figure: 3

**Baseline model:** two people on skis sitting on a snowy surface.
**With object features:** a person standing next to snowboards attached.

Figure: 4



**Baseline model:** a cow is standing in a open field as it grazes.
**With object features:** cows eat alone grazing on grasses in a hill.

Figure: 5

**Baseline model:** man walking next to an old fashioned planes.
**With object features:** a small black and white picture of a prop plane sitting on the runway.

Figure: 6



**Baseline model:** a brown bears perch in front of their mom and another animal.
**With object features:** a brown bear is standing behind a group of brown bears.

Figure: 7

**Baseline model:** two women make homemade my diners can be judged on a table.
**With object features:** a group of people sitting at a blue table of food.

Figure: 8

## CONCLUSIONS

In this paper, we displayed an attention-based Encoder-Decoder picture captioning demonstrates that employments two strategies of highlight extraction, a picture classification CNN (Xception) and a question discovery module (YOLOv4), and demonstrated the adequacy of this plot. We presented the significance figure, which prioritizes frontal area huge objects over foundation little ones, and favours objects with tall certainty over those with moo certainty and illustrated its impact on expanding scores. We appeared how our strategy progressed the scores and compared it to past works within the score increment, particularly the CIDEr metric which expanded by 15.04%, reflecting progressed linguistic saliency.

Not at all like past works, our work recommended to advantage from all protest location highlights extricated from YOLO and appeared the impact of sorting the extracted protest labels. This may be encourage progressed by way better strategies for combining question discovery highlights with convolutional highlights. Future work can also benefit from wealthy protest semantic data from caption writings rather than fair protest formats, which can increment picture captioning precision. Besides, more advanced strategies can be utilized to encode question highlights some time recently contributing them into the decoder, and more complex dialect models, such as Meshed-Memory Transformers can be utilized.

## REFERENCES

1.  Farhadi A, Hejrati M, Sadeghi MA, Young P, Rashtchian C, Hockenmaier J, Forsyth D. Every picture tells a story: generating sentences from images. In: European conference on computer vision. Berlin: Springer; 2010. p. 15–29.

2.  Hochreiter S, Schmidhuber J. Long short-term memory. Neural Comput. 1997;9(8):1735–80.

3.  Cho K, Van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, Bengio Y. Learning phrase representations using RNN encoder-decoder for statistical machine translation. https://arxiv.org/abs/1406.1078. Accessed 3 Jun 2014.

4.  Katiyar S, Borgohain SK. Image captioning using deep stacked LSTMs, contextual word embeddings and data augmentation. https://arxiv.org/abs/2102.11237. Accessed 22 Feb 2021.

5.  Redmon J, Farhadi A. YOLO9000: better, faster, stronger. In: Proceedings of the IEEE conference on computer vision and pattern recognition. Piscataway: IEEE; 2017. p. 7263–71.

6.  Lanzendörfer L, Marcon S, der Maur LA, Pendulum T. YOLO-ing the visual question answering baseline. Austin: The University of Texas at Austin; 2018

7.  Bochkovskiy A, Wang CY, Liao HY. Yolov4: optimal speed and accuracy of object detection. https://arxiv.org/abs/2004.10934. Accessed 23 Apr 2020.

8.  Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. https://arxiv.org/abs/1409.1556. Accessed 4 Sep 2014

9.  Deng J, Dong W, Socher R, Li LJ, Li K, Fei-Fei L. ImageNet: a large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. Piscataway: IEEE; 2009. p. 248–55.

10. Ren S, He K, Girshick R, Sun J. Faster R-CNN: towards real-time object detection with region proposal networks. Adv Neural Inf Process Syst. 2015;28:91–9.