

Improved Image Captioning Using Associative Correlation

M.S.Minu, Akilan Ganesan, Ram Abhilash V, M Mageswaran

Abstract: *The production of appropriate captions to a given image is the characteristic feature of Image captioning. It helps to identify the notable areas of a certain image. Even Though neural networks do have, as of late accomplished promising outcomes, a key issue that still exist is that they can just portray ideas found in trained image-sentence data sets. Proficient learning and portrayal of novel ideas has accordingly been the centre of focus of recent researches. This would help to ease the costly labour of naming data and their sets. The authors propose a focused search by using the data gained from the individual areas of the image with different sources of information from various data sets to train the existing model to come up with captions that are outside the image captioning datasets. Our model utilizes semantic data to produce subtitles for several item classes in ImageNet object identification dataset. Both programmed assessments and human assessments demonstrate that our model significantly outflanks earlier work in having the capacity to portray a lot more classifications of articles which are outside the scope of the existing datasets.*

Keywords: *Image captioning, novel concept, visual attention.*

I. INTRODUCTION

Human cognition plays an important role in learning and identifying innovative ideas. Being a kid, we learnt a lot of concepts, and while growing up, by observing the real world and descriptions given by our parents, we learned to construct appropriate semantic sentences. It has made our language more efficient and understandable [1, 2]. Even though the process was slow and time taking, it has got a lot faster after we have learnt enough concepts and gathered sufficient knowledge. In the field of Image recognition, several methods were proposed [3, 4] to learn using limited examples and build new categories of objects. It is not always possible to come up with novel concepts in a single category and thus it is an important practice to transfer knowledge from categories that have been already learnt. This is because we do not always have sufficient data for creative concepts. The model parameters are extremely huge and, thus, we would not want to re-train the entire model each time again and again just to add a few images.

Manuscript published on 30 April 2019.

* Correspondence Author (s)

Akialn Ganesan*, Department of Computer science, SRM institute of science and technology, Chennai, India.

Ram Abhilash, Department of Computer science, SRM institute of science and technology, Chennai, India.

M Mageswaran, Department of Computer science, SRM institute of science and technology, Chennai, India.

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an [open access](https://creativecommons.org/licenses/by-nc-nd/4.0/) article under the CC-BY-NC-ND license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

The recent years have seen the upcoming of deep recurrent neural network models that show great results in generate captions for images and videos in a proper way . These methods have been successful with the use of very huge corpora of images data sets, like the MSCOCO[5] and Flickr30k[6]. ImageNet[7] describes a smaller variant of objects in relation with the number of objects that can be labelled in object recognition datasets. Even though latest object recognition systems can differentiate hundreds of object classes, the current caption generating models lack the natural language processing part to form semantic sentences[8]. They integrate new sentences whose concepts are previously understood but without any separate examples for image and sentence correlation. This limitation can be overcome by joining visual parts of lexical units by creating a model to generate sentences about image objects which are absent in caption or descriptive sentences (associated image to description) but are present in text (unpaired text data) and object recognition datasets (unpaired image).

Present day visual classifiers[9] can perceive a great many article classifications, few of them are essential or segment based (for example TV), and others which are more detailed and fine-grained (for example dial-telephone, mobile phone). In any case, latest cutting edge image captioning frameworks that gain straightforwardly from pictures and portrayals, depend entirely on matched image description information for supervision and fail in their capacity to sum up and depict this immense arrangement of unmistakable articles in context of the given image. While these frameworks could be improved largely by building larger image/video description datasets, getting such descriptive information would be costly and difficult. Moreover, generating captions is an arduous task since models need to not just accurately recognize visual ideas contained in a picture, yet should likewise make these ideas into a rational sentence.

The concept behind this paper is to use the semantic structure[10] of the given photograph and content to proficiently teach the model novel ideas. For instance, "a deer/okapi in a arid field" shares the setting "in a arid field." Combining zebra and giraffe with the proper balancing setting is legitimately right. Such a semantic structure empowers a progressively effective approach to teach a novel idea to the model. Assume the framework has taken in the idea of a deer however has never observed a okapi, it can figure out how to portray a okapi in a arid field, just by perceiving deer and knowing the way that a okapi can be in a arid field like a deer.



The image and caption are in this manner are unpaired—the required information of novel ideas comprise of: a free image base giving visual data and an autonomous learning base providing logic information.

The proposed model uses the recent deep captioning models while combining it with conventional and recurrent networks for better semantic captioning. The existing models could only caption the objects found in the image and descriptive sentence data-set, and it is revolutionary that the our framework can generate sequences of sentences or captions about the new or novel image objects by fusing them with already seen expressions.

To describe the image more efficiently, the newly found objects in our model follow two methodologies. First, it contains of its own pair of visual classifier (CNN) and Sentence Constructor (RNN). These can independently be trained on image data-sets that are unpaired along with unpaired language data. Moreover, the image classifier as well as the Sentence Constructor can be integrated into a deep machine learning model which can primarily be used to train paired image-caption data efficiently. The later part is important as it uses multi-modal layer for generating captions with novel objects. The required objects are having an associative correlation between image and sentences datasets. It can be sent over to novel objects exclusively available in unpaired image and language datasets. In this paper, we use external text to create relations with new objects with concepts available in associated data and propose other options to transfer correlated data from known objects to novel objects.

The another enormous test accompanies the substitution of a picture of any novel item. In contrast to supplanting words, which should be possible with high accuracy, supplanting districts in a picture is troublesome because of the unacceptable exactness of acknowledgement and division. That is, we can't actually expel the visual parts of pre existing ideas and fill in the visual data or vectors of novel ideas on the pixel level. Hence the work on abnormal state semantic highlights rather than pixels we make another article class which goes about as a pre existing for future referencing of that novel item.

II. RELATED WORKS

A. Image Captioning:

In the last decade or so there has been profound advancement in neural systems for common language and machine acknowledgement capabilities. For characteristicgge, Recurrent Neural Networks(RNNs) and Long Short Term Memories(LSTMs) accomplish the cutting edge execution for some, NLP tasks, for example, machine interpretation and discourse acknowledgement. For Computer vision, profound Convolutional Neural Networks (CNNs) outflank past strategies by an extensive edge for the assignments of item grouping and discovery. As of late an assortment of CNN frameworks have achieved extraordinary results on the image writing undertaking. They pursue a CNN - RNN framework[11]: the process begin as abnormal state highlights are extricated from the CNN prepared on the picture order errand, and after that an intermittent model figures out how to anticipate ensuing expressions of an

inscription adapted on picture highlights and recently anticipated words. Our model is like such CNN-RNN structures regardless, neither of these models can be prepared to identify new objects that aren't present already in the available dataset.

B. Multi-modal Learning:

The techniques for picture sentence retrieval, picture portrayal age and visual inquiry noting have grown quick as of late. Many of them adopt an CNN-RNN the structure that advances the log-probability of the caption given the picture, and trains the systems in a start to finish the way. In special cases are models which embrace a multi-modal machine learning model[12] where intermittent language highlights and picture highlights are inserted in a multi-modal space which wires visual discoverers, language models, and multi-modal resemblance models in a high-efficiency pipeline. The multi-modal implanting is then utilized to foresee the description by each word. Despite the fact that our methodology utilizes distributional word embeddings, our model contrasts as in it tends to be prepared with unpaired content and visual information yet at the same time consolidate the semantic data at a later stage amid caption generation.

C. Zero Shot Learning:

In zero-shot learning, image features are associated with dense vectors of words. The thick word vectors in zero-shot learning are pre-arranged from a great deal of web corpus dataset and the word dataset depiction is gotten from co-occasion with various words. Learning new objects is possible from only a selected range of examples[4, 13]. Be that as it may, these work just consider words or traits rather than sentences still it has received substantial attention in computer vision since it winds up hard to get adequate checked pictures as the amount of article classes creates. In particular, our strategy is derived from the past zero-shot learning work that basically searches for object connections in outer content corpora to decide how protests are identified with one another, at that point the obscure items and known objects are classified based on their relationship.

III. PROPOSED MODEL

We allude to the ideas in the first matches as pre existing ideas and the ideas to be gained from the unpaired information as new or novel ideas. The language generating grammar in picture subtitling is simple, and the pre existing sentences are sufficient to cover the syntax. When learning new or novel ideas, the greater part of the exertion is on figuring out how to remember them, not adapting new language or linguistic structure. The proposed model produces novel sentences about items concealed in combined picture sentence information. Despite the fact that usually to pre-train profound inscription models on unpaired picture information, in contrast to existing models, we can depict objects present in unpaired picture information however not present in matched picture sentence information. The proposed model disintegrates into 3 pre-prepared models viz. Visual classifier (CNNs), Joint model (CNN+RNN), NLP (LSTMs)[17]. The display begins by removing object classes from the picture utilizing a visual classifier (CNNs), here the

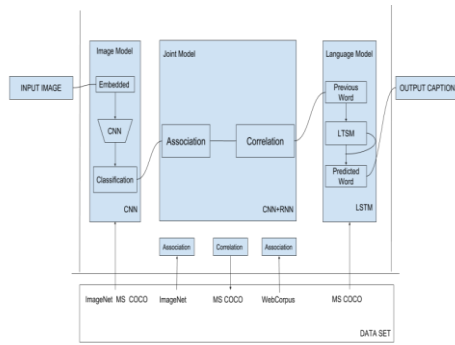


Fig.1:Architecture

image and object data is given to a joint model consisting of an image-class data and a text corpus. The joint model is zero shot trained so that the text corpus can associate with a specific object class from image-class data, this happens only if the object class is present in image-class data but absent in image-sentence data. As this is a novel data to the image sentence data the association from the joint model is used to modify the sentence from the pre existing paired image-sentence data to generate visual features and sentence pair for the novel object i.e correlate the novel object to associated objects visual features to generate a new class in the image-sentence pairs. The caption is then generated using the object classes given by the joint model using a NLP (LSTMs).

A. Image Model:

Our visual acknowledgement display is a neural system and is prepared on article acknowledgement datasets. Not at all like normal visual acknowledgement models that are prepared with a solitary mark on a grouping errand, for the assignment of picture inscribing a picture demonstrate that has high certainty over numerous visual ideas. It is desirable over have the picture occurrence at the same time. Consequently, we train our model utilizing various names with a multi-mark misfortune.

B. Language Model:

Our language framework completely depends on LSTM[14](i.e Long Short Term Memory) as well as the RNNS(i.e Recurrent Neural Networks). The language display is prepared to foresee the following word w_t in a given arrangement of words w_0, \dots, w_{t-1} .

3.3. Join Model:

The objective of the joint model is to produce a sentence conditioned on an image which contains a novel object. The joint model can be broken down into 2 major phases association and correlation. Association refers to identifying of a novel object and associating it with a pre existing object class in ImageNet[6]. And Correlation refer to creation of the novel object class with necessary visual features and sentence pairs for future captioning of the novel object

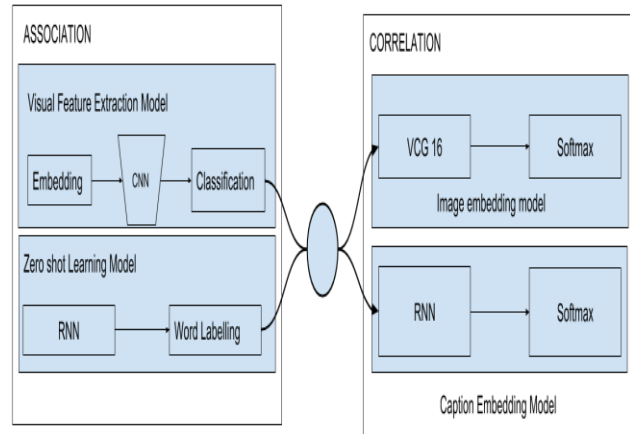


Fig.2:Associative Correlation

C. Association Phase:

To ensure the accuracy of the proposed model, the generation of novel object classes ought to be limited to a sensible range. We propose a strategy for estimating object comparability dependent on the reference information. The substantial number of conceivable novel objects are filtered using the similarities. The filtering can be changed to influence the rate at which the model generates novel object classes. This creates a trade off between accuracy and novel class generation. Visual feature extraction:- To extract visual features we encode associated images from ImageNet. We segment the image into area or blocks, and encode each block using VGG-16[9] to acquire a lot of highlight vectors that store neighbourhood visual data

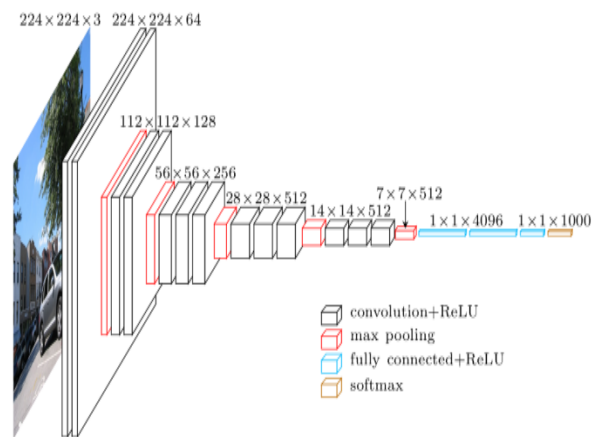


Fig.3:Association Mod

D. Correlation Phase:

To ensure the capability of the proposed model to generate a caption for a novel object without going through the association phase again, the generation of image-sentence pairs is required. This includes sentence generation and image generation. Sentence generation:-An innocent path is to supplant the already defined object name by the novel object name in every sentence of which contains already defined object name. Be that as it may, such a savage power approach almost never succeeds to take in the specific setting of the sentence; as such, it might incite freakish sentences, for example, "A sheep is being sheared in a shed" → "A puppy is being sheared in a shed". With the thought of explicit settings in sentence age applicants, we perform exact substitution to create increasingly sensible sentences. Picture age:- Unlike supplanting words in a sentence, supplanting ideas in a picture are progressively muddled. We accomplish this by picture mixing i.e expelling out the visual highlights of the already defined idea from the picture and after that adding in the novel idea[19]. The exact limits of the already defined idea is resolved utilizing striking sifting.

Algorithm 1 Context-Specific Replacement Strategy

Input: cn: novel concept
 C: set of seed concepts
 KB: knowledge base
 Ss: set of sentences in human-labeled data
 Nmax: maximum number of replacements

1. Find the top 5 similar concepts to cn in C, and denote as Cs
2. Initialize $S_n \leftarrow \emptyset$
3. for S in Ss do
4. if S contains any concept cs in Cs then
5. do POS-tagging for S, obtaining (cs, b)
6. if (cn, b) \in KB then
7. replace words about cs with words about cn in S, obtaining the synthetic sentence S
8. $S_n \leftarrow S_n \cup \{S\}$
9. if $|S_n| > N_{max}$ then
10. end loop
11. **return** S_n : the set of pseudosentences of cn

Fig.4:Algorithm

IV. DATASETS

We use ImageNet[6], MSCOCO[5] and an external text corpus[18] for the training of the proposed model. The ImageNet is used as the object detection dataset. As it contains a whopping 14,197,122 images which has more than 21841 object classes. This improves the chance of association for novel objects. The MS-COCO i.e common object in context dataset is used as the main captioning image-sentence pair dataset. It consists of 123,287 images

with 5 captioning reference sentences for each picture. And the external text corpus is created by using a WebCropus tool on major sites influential to the vocabulary of the language model and word-based association.

A. External Text Corpus (Web Corpus):

We extricated sentences from Wikipedia, Gigaword, and the British National Corpus (BNC) utilizing [16]. This dataset was utilized to prepare the LSTM language display. For the thick word portrayal in the proposed model for affiliation, we use GloVe[15] pre-prepared on outside corpora including Gigaword and Wikipedia. To make our language display vocabulary we recognized the 80k most incessant tokens from the consolidated outer corpora. We refine this vocabulary further to a lot of 72,700 words that likewise had GloVe embeddings. This refinement is important to have expected sentences to relate after a word is related utilizing GloVe.

B. Image Caption data:

This dataset for the paired image and sentence data is taken from MSCOCO by bunching the essential 80+ object characterizations using the cosine separator on word to vector function (of article mark) and picking one thing from each gathering to hold out from getting ready.

C. Image data:

The picture informational collection is perceived as articles that can be spoken to in both the ImageNet and our word or object set (vocabulary), however not in MSCOCO. The words chose have an assortment of sets of classes which contain detailed and fine grained characterization e.g., "dog puppy" and "chrysanthemum", descriptors e.g., "chiffon", "Wool-en", and section level words for e.g., "rodent"..

V. EXPERIMENTATION

To exhibit the capacity of the proposed model, we depict questions from the ImageNet[6] dataset where no matched picture sentence information exists. New Objects: Many Objects that are available in ImageNet just as the model's object identification vocabulary have yet to be referenced in the image-caption dataset i.e MSCOCO[5]. For quantitative assessment we measure the level of articles for which the model can portray no less than one picture of the item. Besides, we get human assessments contrasting our model and past work on whether the model can fuse the item mark genuinely in the portrayal together with how well it depicts the picture

A. Captioning Novel Object:

In spite of the fact that we utilize the equivalent CNN-RNN framework[11], the affiliated relationship model is both prepared to delineate more classes and precisely coordinate new words into depictions even more as frequently as could reasonably be expected. Information picture sets can bomb either in respect to find a sensible article that is semantically and phonetically like the novel or new thing, or concerning the language making a sentence using the thing name, in our model the former never occurs (for instance it's not required to unequivocally perceive near things), diminishing the general wellsprings of mistake.

B. Human Assessment:

ImageNet pictures don't have going with inscriptions and this makes the assignment significantly more difficult to assess. To think about the execution of our proposed model we acquire human decisions on inscriptions produced by the models on a few item classifications. While choosing the pictures, for article classes that both (proposed and more seasoned) models can depict, we make a point to choose no less than two pictures for which the two models notice the item name in the portrayal. Each picture is displayed to three specialists. We directed two human investigations (test interface is based on the enhancement): Given the picture, the ground-truth image object characterization (and semantics), and the descriptive sentences made by our associative-correlation model, we evaluate the captioning based on:

C. Word Consolidation:

We require the individuals to pick any sentence/inscription which fuses the item mark genuinely in the depiction. The choices gave are: (i) Sentence 1 choice of word is better compared to 2, (ii) Sentence 2 choice of word is better compared to 1, (iii) Both sentences consolidate the word similarly well, or (iv) Neither of them did well.

D. Visual Description:

Humans can be utilized to choose a superior sentence among both which characterizes the image better. This empowers us to take a gander at both how efficiently the model combines the novel or new thing in the sentence, similarly as how relevant the portrayal is to the picture. On the given list of pictures comparing to the image objects that the two models can depict (Intersection), both show up equally coordinated. In any case, seeing all article classifications (Union), the proposed model can both circuit the thing mark in the caption, and portray the image more defined and accurately than past models.

Table. 1

Objects subset →	Word Incorporation		Image Description	
	Union	Intersection	Union	Intersection
NOC is better	43.78	34.61	59.84	51.04
DCC is better	25.74	34.12	40.16	48.96
Both equally good	6.10	9.35	-	-
Neither is good	24.37	21.91	-	-

VI. CONCLUSION

In our paper, we acquainted the model to distinguish and utilize novel items in picture subtitling. It uses data from unpaired information and reuses punctuation in the first combined information. With a cooperative relationship, the subtitling model can effectively learn novel ideas with no human-marked picture sentence pair, in this manner, significantly diminishing the expense of expanding the idea go. The unpaired information is made up from an outer content corpus base and a huge picture base. Then We have developed the content based web slithering from the significant web page and separating them. The picture data set (ImageNet) is sorted out as classifications to such an extent that it is conceivable to question an idea to bring the comparing pictures. The trial results obtained when conducted on a constrained list of MSCOCO images demonstrated that the proposed methodology gives critical enhancements over cutting edge techniques as far as distinguishing new articles and language quality. Our methodology is strong to the measure of demonstrating stable execution inside a variably large scope of hyper parameters.

REFERENCES

1. D. Swingley. Fast mapping and slow mapping in children's word learning. *Language learning and Development*, 6(3):179–183, 2010
2. T. H. Heibeck and E. M. Markman. Word learning in children: An examination of fast mapping. *Child development*, pages 1021–1034, 1987..
3. T. Tommasi, F. Orabona, and B. Caputo. Learning categories from few examples with multi model knowledge transfer. *TPAMI*, 36(5):928–941, 2014.
4. A. Lazaridou, E. Bruni, and M. Baroni. Cross-modal mapping between distributional semantics and the visual world. In *ACL*, pages 1403–1414, 2014.
5. T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014.
6. J. Wang, "Fundamentals of erbium-doped fiber amplifiers arrays (Periodical style—Submitted for publication)," *IEEE J. Quantum Electron.*, submitted for publication.
7. B. Plummer, L. Wang, C. Cervantes, J. Caicedo, J. Hockenmaier, and S. Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.
8. L. A. Hendricks, S. Venugopalan, M. Rohrbach, R. Mooney, K. Saenko, and T. Darrell. Deep compositional captioning: Describing novel object categories without paired training data. *arXiv preprint arXiv:1511.05284*, 2015
9. K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
10. A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3128–3137, 2015.
11. O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. *CVPR*, 2015
12. R. Kiros, R. Salakhutdinov, and R. S. Zemel. Unifying visual-semantic embeddings with multimodal neural language models. *TACL*, 2015.
13. S. Antol, C. L. Zitnick, and D. Parikh. Zero-shot learning via visual abstraction. In *ECCV*, pages 401–416. 2014.

Improved Image Captioning Using Associative Correlation

14. Huang, X., Tan, H., Lin, G., & Tian, Y. A LSTM-based bidirectional translation model for optimizing rare words and terminologies. 2018 International Conference on Artificial Intelligence and Big Data (ICAIBD), 185-189. 2018.
15. Pennington, J., Socher, R., & Manning, C.D. Glove: Global Vectors for Word Representation. EMNLP. 2014
16. Manning, C.D., Surdeanu, M., Bauer, J., Finkel, J.R., Bethard, S., & McClosky, D. The Stanford CoreNLP Natural Language Processing Toolkit. ACL. 2014.
17. Sepp Hochreiter and Jurgen Schmidhuber. Long short-term " memory. Neural Comput., 1997.
18. R. Mihalcea, C. Corley, C. Strapparava et al., "Corpus-based and knowledge-based measures of text semantic similarity," in AAAI, vol. 6, pp.775-780, 2006.
19. Ryan Kiros, Ruslan Salakhutdinov, and Richard S. Zemel, "Unifying visual-semantic embeddings with multimodal neural language models," CoRR, vol. abs/1411.2539, 2014.