# Dynamic Caption Generation with Image Annotation

Thirupathi Podeti[1], Bhanu Prasad A[2]

[1]*M.Tech Student, Dept. of CSE, Vardhaman College of Eng., Hyderabad, India*
[2]*Associate Professor, Dept. of IT,Vardhaman College of Eng., Hyderabad, India*

**Abstract - Now a day's, whenever retrieving footage from the search Engines that retrieves footage whereas not analyzing their content, simply by matching user queries against the image's file name and format, user-annotated tags, captions, and, generally, text shut the image. To boot the retrieved image does not contain any matter data beside the photographs. We have a bent to introduced the task of automatic caption generation for news footage. The task fuses insights from computer vision and communication method and holds promise for various multimedia applications, like image retrieval, development of tools supporting fourth estate management, and for folks with handicap. It's potential to search out a caption generation model from frail labelled data whereas not costly manual involvement. Instead of manually creating annotations, image captions unit of measurement treated as labels for the image. Although the caption words unit of measurement avowedly noisy compared to ancient human-created keywords, we have a bent to indicate that they're going to be used to learn the correspondences between visual and matter modalities, and to boot operate a gold traditional for the caption generation task. We have given extractive and hypothetical caption generation models. A key aspect of our approach is to allow every the visual and matter modalities to influence the generation task.**

**Keywords - Captions, Annotation, Multimedia Applications.**

## I. INTRODUCTION

Now a day's, whenever retrieving footage from the search Engines that retrieves footage whereas not analyzing their content, just by matching user queries against the image's file name and format, user-annotated tags, captions, and, generally, text shut the image. Also the retrieved image doesn't contain any matter information beside the pictures. We have a bent to introduced the task of automatic caption generation for news footage. The task fuses insights from pc vision and communication technique and holds promise for numerous transmission applications, like image retrieval, development of tools supporting fourth estate management, and for people with handicap. it's potential to go looking out a caption generation model from frail tagged information whereas not expensive manual involvement. rather than manually making annotations, image captions unit of activity treated as labels for the image. though the caption words unit of activity avowedly clattery compared to ancient human-created keywords, we have a bent to point that they are aiming to be accustomed learn the correspondences between visual and matter modalities, and also operate a gold ancient for the caption generation task. We have given extractive and theoretic caption generation models. A key facet of our approach is to permit each the visual and matter modalities to influence the generation task.

## II. RELATED WORK

In this paper, we tackle the related problem of generating captions for news images. Our approach leverages the vast resource of pictures available on the web and the fact that many of them naturally co-occur with topically related documents and are captioned. We focus on captioned images embedded in news articles, and learn both models of content selection and surface realization from data without requiring expensive manual annotation. At training time, our models learn from images, their captions, and associated documents, while at test time they are given an image and the document it is embedded in and generate a caption. Compared to most work on image description generation, our approach is shallower, it does not rely on dictionaries specifying image-to-text correspondences, nor does it use a human-authored grammar for the caption creation task. It uses the document co-located with the image as a proxy for linguistic, visual, and world-knowledge. Our innovation is to exploit this implicit information and treat the surrounding document and caption words as labels for the image, thus reducing the need for human supervision. During testing, we are given a document and an associated image for which we must generate a caption. And the knowledge base must contain two types of information, information about how the images (or image regions) corresponds to words and information about how these words can be combined to create a human-readable sentence.

## III. LITERATURE SURVEY

A. Image Classification for Content-Based Indexing

User queries in content-based retrieval are typically based on semantics and not on low-level image features. Providing high-level semantic indices for large databases is a challenging problem. We have shown that certain high-level semantic categories can be learnt using specific low-level image features under the constraint that the images do belong to one of the classes under consideration

B. Content-Based Image Retrieval at the End of the Early Years

The paper presents a review of 200 references in content-based image retrieval. The paper starts with discussing the working conditions of content-based retrieval: patterns of use, types of pictures, the role of semantics, and the sensory gap. Subsequent sections discuss computational

steps for image retrieval systems. Step one of the review is image processing for retrieval sorted by color, texture, and local geometry. Features for retrieval are discussed next, sorted by: accumulative and global features, salient points, object and shape features, signs, and structural combinations.

### C. Efficient Object Recognition Using Boundary Representation and Wavelet Neural Network

In this paper, an efficient object recognition method using boundary representation and the wavelet neural network is proposed. The method employs a wavelet neural network (WNN) to characterize the singularities of the object curvature representation and to perform the object classification at the same time and in an automatic way.

### D. Simultaneous Image Classification and Annotation

In this paper, we develop a new probabilistic model for jointly modeling the image, its class label, and its annotations. Our model treats the class label as a global description of the image, and treats annotation terms as local descriptions of parts of the image. Its underlying probabilistic assumptions naturally integrate these two sources of information.

### E. I2T: Image Parsing to Text Description

In this paper, we present an image parsing to text description (I2T) framework that generates text descriptions of image and video content based on image understanding. The proposed I2T framework follows three steps: 1) input images(or video frames) are decomposed into their constituent visual patterns by an image parsing engine, in a spirit similar to parsing sentences in natural language; 2) the image parsing results are converted into semantic representation in the form of Web ontology language (OWL).

### F. Baby Talk: Understanding and Generating Simple Image Descriptions

This paper exploits both of these lines of attack to build an effective system for automatically generating natural language – sentences – from images. It is subtle, but several factors distinguish the task of taking images as input and generating sentences from tasks in many current computer vision efforts on object and scene recognition.

## IV. IMPLEMENTATION DETAILS

### A. Data Collection

We created our own dataset by downloading articles from the News websites. The dataset covers a wide range of topics including national and international politics, technology, sports, education, and so on. News articles normally use color images which are around 200 pixels wide and 150 pixels high. The captions tend to use half as many words as the document sentences and more than 50 percent of the time contain words that are not attested in the document.
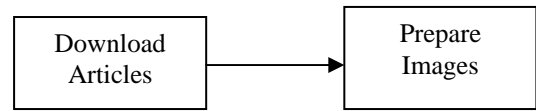


Fig .1 Data Collection

### B. Input Preparation

The document should contain the necessary background information which the image describes or supplements. And also we can exploit the rich linguistic information inherent in the text and address caption generation with methods relative to text summarization without extensive knowledge engineering. The caption generation task is not constrained in any way, words and syntactic structures are chosen with the aim of creating a good caption rather than rendering the task acceptable to current vision and language generation techniques.
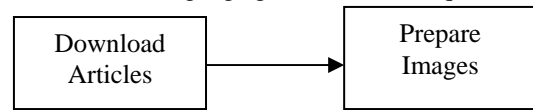


Fig .2 Input Preparation

### C. Abstractive Caption

We turn to abstractive caption generation and present models based on single words but also phrases. Content selection is modeled as the probability of a word appearing in the headline given that the same word appears in the corresponding document and is independent of other words in the headline. They also take the distribution of the length of the headlines into account in an attempt to relative to the model toward generating output of reasonable length.
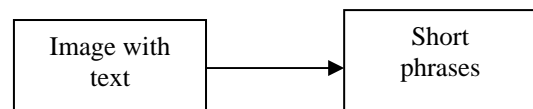


Fig .3 Abstractive Caption

### D. Extractive Caption

This Extractive caption mostly focuses on sentence extraction. The idea is to create a summary simply by identifying and subsequently concatenating the most important sentences in a document. Without a great deal of linguistic analysis, it is possible to create summaries for a wide range of documents, independently of style, text type, and subject matter. For our caption generation task, we need only extract a single sentence. And our guiding hypothesis is that this sentence must be maximally similar to the description keywords generated by the annotation model.
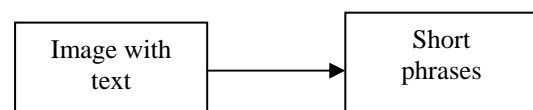


Fig. 4 Extractive Caption

## V.  RELATED WORK

Search engines deployed on the net retrieve pictures while not analyzing their content, just by matching user queries against collocated matter data. Examples embrace information (e.g., the image's file name and format), user-annotated tags, captions, and, generally, text encompassing the image. As this limits the relevance of search engines (images that don't coincide with matter knowledge cannot be retrieved), an excellent deal of labor has targeted on the event of ways that generate description words for an image mechanically. This approach will produce sentences of top quality that area unit each important and fluent; but, the reliance on manually created resources mostly limits the readying of existing ways to real-world applications. Developing dictionaries that specify thoroughly image-to-text correspondences may be a tough and long task that has got to be recurrent for brand spanking new domains and languages. Constant is additionally true for templates and rule-based generation methods; the previous area unit generally specific to the domain in question and not transportable to new tasks, whereas the latter are often additional general and lingual subtle, albeit with in depth data engineering.

Extractive and Abstractive Caption

Extractive caption generator attracts inspiration from previous work on automatic report, most of that focuses on sentence extraction. The concept is to form a outline just by distinctive and after concatenating the foremost necessary sentences during a document. While not a good deal of linguistic analysis, it's attainable to form summaries for a good vary of documents, severally of fashion, text type, and subject material. For our caption generation task, we want solely extract one sentence. And our guiding hypothesis is that this sentence should be maximally almost like the outline keywords generated by the annotation model. The 'phased' nature of extractive activities implies that several, or all, phases of AN extractives operation –exploration, evaluation, development, production and decommissioning - will occur at identical time. Production can habitually start whereas development continues. Exploration and analysis designed to spot and expand reserves (or to 'prove up' reserves) can continue. Major expansionary expenditure is usually postponed to maximise money flows (e.g., extending a decline or shaft, up-scaling plant capability, increasing well heads), solely being completed as needed throughout the lifetime of the operation.

The design of the mining or oil and gas operation seeks to maximise returns from the operation overall and doesn't expressly think about individual 'expansions' as production continues, e.g., during a mining operation, ore could also be wasteful on its own however 'blended' with different ore to realize optimum ore grades for process, or otherwise wasteful ore in AN underground mine, could also be strip-mined in gaining access to high grade ore; and therefore the initial infrastructure is built in sight of the ultimate scale of the operation, and so it's going to not be supported by money flows expected from accessible reserves which may presently be extracted.

## VI.  RESULT

In this, we explored the possibility of automatically generating captions for news Images. This software allows both visual and textual information to influence automatic caption generation task. Here, image annotation model converts features of an image into self explanatory keywords which are subsequently used to guide the generation process. Given an accompanying document our abstractive model generates a sentence.

ARTICLE:

The British government could face claims it violated international humanitarian laws by allowing US arms flights to Israel to use UK airports. The Islamic Human Rights Commission (IHRC) is seeking permission to contest government bodies over what it says are crimes against the Geneva Convention. A number of US planes said to be carrying bombs to Israel refuelled in the UK during the Lebanon conflict. The IHRC said it received complaints from Britons with families in Lebanon. The commission is accusing the government of "grave and serious violations" of international humanitarian law. It is seeking permission to bring its case against the Civil Aviation Authority, the Foreign and Commonwealth Office and Defense Secretary Des Brown in the High Court. The IHRC said it is bringing the case after receiving "many complaints ... from British citizens whose family members are in Lebanon and facing grave danger as well as acts of terror". The US aircraft believed to have refueled in the UK are said to have been carrying supplies including "bunker buster bombs".

Fig. 5 Example Image for Article

## VII.  CONCLUSION

We have introduced the task of automatic caption generation for news pictures. The task fuses insights from laptop vision and language process. Alike from previous work, we have approached this task during a knowledge-lean fashion by victimization the huge resource of pictures offered on the web and exploiting the very fact that a lot of of those go with matter data (i.e., captions and associated documents). Our results show that it's attainable to find out a caption generation model from sapless tagged knowledge while not expensive manual involvement. Image captions area unit treated as labels for the image. though the caption words area unit true screaky compared to ancient human-created keywords, we have a tendency to show that they will be accustomed learn the correspondences between visual and matter modalities, and additionally function a gold commonplace for the caption generation task.

Moreover, this news dataset contains a novel part, the news document, that provides each data relating to to the image's content and wealthy linguistic data needed for the generation procedure.

## VIII.    FUTURE WORK

As video process typically involves process key frames (images) from streaming video knowledge, it's conjointly attainable to adapt existing models and applications from pictures to video (e.g., automatic video summarization).The dataset mentioned during this thesis will be additional refined in step with specific Applications to use the news document to extend the annotation keywords by distinctive synonyms or perhaps sentences that ar just like the image caption. a comprehensible extension would be taking abstraction data into consideration once managing image representations. Currently, we tend to treat the image regions or detected regions of interest as bags-of-words, that may be extended to bigrams in step with their abstraction relations. As an example, we tend to might experiment with options associated with document structure like titles, headings, and sections of articles and conjointly exploit syntactical data additional directly. We could, however, improve grammaticality additional globally by generating a grammatical tree (or dependency graph).

### REFERENCES

[1]   A. Vailaya, M. Figueiredo, A. Jain, and H. Zhang, "*Image Classification for Content-Based Indexing*," IEEE Trans. Image Processing, vol. 10, no. 1, pp. 117-130, 2001.

[2]   A.W. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, "*Content-Based Image Retrieval at the End of the Early Years*,"

IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 22, no. 12,pp. 1349-1380, Dec. 2000.

[3]   P. Duygulu, K. Barnard, J. de Freitas, and D. Forsyth, "*Object Recognition as Machine Translation: Learning a Lexicon for a Fixed Image Vocabulary*," Proc. Seventh European Conf. Computer Vision, pp. 97-112, 2002.

[4]   D. Blei, "*Probabilistic Models of Text and Images*," PhD dissertation, Univ. of Massachusetts, Amherst, Sept. 2004.

[5]   K. Barnard, P. Duygulu, D. Forsyth, N. de Freitas, D. Blei, and M. Jordan, "*Matching Words and Pictures*," J. Machine Learning Research, vol. 3, pp. 1107-1135, 2002.

[6]   C. Wang, D. Blei, and L. Fei-Fei, "*Simultaneous Image Classification and Annotation*," Proc. IEEE Conf. Computer Vision and Pattern Recognition, pp. 1903-1910, 2009.

[7]   V. Lavrenko, R. Manmatha, and J. Jeon, "*A Model for Learning the Semantics of Pictures*," Proc. 16th Conf. Advances in Neural Information Processing Systems, 2003.

[8]   S. Feng, V. Lavrenko, and R. Manmatha, "*Multiple Bernoulli Relevance Models for Image and Video Annotation*," Proc. IEEE Conf. Computer Vision and Pattern Recognition, pp. 1002-1009, 2004.

[9]   L. Ferres, A. Parush, S. Roberts, and G. Lindgaard, "*Helping People with Visual Impairments Gain Access to Graphical Information through Natural Language: The igraph System*," Proc. 11th Int'l Conf. Computers Helping People with Special Needs, pp. 1122-1130, 2006.

[10]   A. Abella, J.R. Kender, and J. Starren, "*Description Generation of Abnormal Densities Found in Radiographs*," Proc. Symp. Computer Applications in Medical Care, Am. Medical Informatics Assoc., pp. 542-546, 1995.

[11]   A. Kojima, T. Tamura, and K. Fukunaga, "*Natural Language Description of Human Activities from Video Images Based on Concept Hierarchy of Actions*," Int'l J. Computer Vision, vol. 50, no. 2, pp. 171-184, 2002.