



IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 9 Issue: VII Month of publication: July 2021

DOI: https://doi.org/10.22214/ijraset.2021.36815

www.ijraset.com

Call: 🕥 08813907089 🔰 E-mail ID: ijraset@gmail.com



# **Assistive Vision Technology using Deep Learning Techniques**

Dr. Neeta Verma<sup>1</sup>, Tushar Kumar<sup>2</sup>, Shivam Upadhayay<sup>3</sup>, Tushar Chaudhary<sup>4</sup>, Shivam Singh Chauhan<sup>5</sup> <sup>1, 2, 3, 4, 5</sup>Dept. of Computer Science and Engineering, Inderprastha Engineering College, Ghaziabad, India

Abstract: One of the most important functions of the human visual system is automatic captioning. Caption generation is one of the more interesting and focused areas of AI, with numerous challenges to overcome. If there is an application that automatically captions the scenes in which a person is present and converts the caption into a clear message, people will benefit from it in a variety of ways. In this, we offer a deep learning model that detects things or features in images automatically, produces descriptions for the images, and transforms the descriptions to audio for louder readout. The model uses pre-trained CNN and LSTM models to perform the task of extracting objects or features to get the captions. In our model, first task is to detect objects within the image using pre trained Mobilenet model of CNN (Convolutional Neural Networks) and therefore the other is to caption the pictures based on the detected objects by using LSTM (Long Short Term Memory) and convert caption into speech to read out louder to the person by using SpeechSynthesisUtterance interface of the Web Speech API. The interface of the model is developed using NodeJS as a backend for the web page. Caption generation entails a number of complex steps, including selecting the dataset, training the model, validating the model, creating pre-trained models to check the images, detecting the images, and finally generating captions.

Keywords: CNN, LSTM, Speech Synthesis Utterance, Image caption generator

#### I. INTRODUCTION

Visual impairment, also known as vision defacement or vision loss, is a loss of ability to determine to the point where it causes problems that cannot be corrected by conventional means such as glasses. There are roughly 285 million visually impaired persons in the globe, with over 39 million blind people, according to the World Health Organization[1]. Living with a visual defect can be difficult because many everyday situations are diffilt to understand without good acuity..

Automatically describing the content of images using natural languages can be a difficult and time-consuming task. It has a huge potential impact, for example, it could help people who are blind or visually challenged have a better understanding of their surroundings. It could also provide more accurate and concise descriptions of their surroundings. This project achieves this challenge using deep learning neural networks. The method generates semantically meaningful and grammatically correct captions for the images using information from picture and caption pairings.



Figure 1: Image caption generation using deep learning

We are using deep learning to support the thought of performing on the Image Caption Generator. Automatically describing the contents of an image with suitable English phrases is a difficult challenge, but it has the potential to greatly assist visually impaired people in recognising and acknowledging their environment. The latest mobile phones or cameras can record photographs of the surroundings making it possible for the dim sighted people to form images of their environments. These images will be used to generate captions that will be read aloud to the visually impaired in order that they will get a far better sense of what's happening around them. In this paper, We present a deep recurrent architecture that creates brief descriptions of incoming images automatically. Our model uses a pre-trained model named MobileNet to detect the objects from the image, which contains the functionality of the convolutional neural network (CNN). These features are then fed into a Long Short Term Memory (LSTM) network to get an outline of the image invalid English and SpeechSynthesisUtterance interface of the online Speech API to convert caption into speech to read out louder to the person. The interface of the model is developed using HTML, CSS and JavaScript for the frontend and NodeJS as a backend for the online page. Our model achieves state-of-the-art performance and generates highly explanatory and descriptive captions which will potentially help the needful and improve their living.



ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.429 Volume 9 Issue VII July 2021- Available at www.ijraset.com

#### II. LITERATURE SURVEY

Most work in visual recognition has originally focused on image classification, i.e., assigning labels and features corresponding to a fixed number of categories to images. Great progress in image classification has been made over the past years especially with the use of deep learning techniques[2][3]. Despite the fact that a category label only provides limited information about an image, more detailed descriptions will benefit visually impaired people. Farhadi et al. and Kulkarni et al. [4][5] have attempted to generate more complex image descriptions in the past, but their algorithms rely on hard-coded phrases and visual concepts. Furthermore, the majority of these works strive to accurately convey the contents of an image in a single word. However, this one sentence limits the quality of the descriptions. Several studies, such as those by Li et al., Gould et al., and Fidler et al. [6]-[9], have focused on gaining a comprehensive knowledge of scenes and objects represented in photographs. The aim of these works was to correctly generate the labels corresponding to a fixed number of categories to the type of an image, instead of generating higher-level explanations of the scenery and entity depicted on an image.

Computer vision and natural language processing are linked by a fundamental problem in artificial intelligence that automatically describes the content of an image. First and foremost, methods extract annotations (nouns and adjectives) from images[10][11] then make a sentence out of the annotations[12][13] created a large-scale recurrent convolutional architecture for visual learning and demonstrated the models' utility on three different tasks: video identification, image description, and video description. These models incorporate long-term dependencies into network status updates and are trainable from beginning to end. The difficulty in comprehending the intermediate result is a limitation.

Several works attempt to solve this task by finding the image in the training set that is most similar to the test image and then returning the caption associated with the test image Jia et al., Kuznetsova et al., and Li et al. find multiple similar images, and combine their captions to generate the resulting caption. The LRCN method is further refined for video caption generation [14]-[16]. Vinyals et al. [15] suggested a neural image caption model that is exclusively used for image caption generation, rather than one architecture for three tasks in LRCN. This model is trained to increase the likelihood of the goal description sentence given the training photos by combining GoogLeNet with a single layer of LSTM. The model's performance is assessed both qualitatively and quantitatively. The MS COCO Captioning Challenge (2015), in which the results were judged by humans, ranked this method first. When LRCN and NIC are compared, three differences emerge that may indicate performance differences.

Kuznetsova et al. and Gupta et al. combined object detection and feature learning with a fixed sentence template [5][17][18]. They attempted to identify objects and features in the image, and then used their sentence template to create sentences describing the image based on the identified objects. Nonetheless, this method severely limits the model's output variety.

There has been a renaissance of interest in picture caption generation as a result of recent deep learning advancements[2] [19]-[22]. Several deep learning algorithms have been developed for producing higher-level word descriptions of images [21][22]. Convolutional Neural Networks (CNNs) have been shown to be useful models for tasks such as picture categorization and object detection. New models for obtaining low-dimensional vector representations of words, such as word2vec and GloVe (Global Vectors for Word Representation), as well as Recurrent Neural Networks, Karpathy and colleagues discovered that image features and language modelling can be combined to create models that generate image descriptions.



Figure 2: Extracting the features from the image using CNN.

Christopher Elamri, Teun de Planque [23] also demonstrates a deep recurrent architecture that automatically generates succinct explanations for photos. To extract features from images, our models use a convolutional neural network (CNN). These attributes are then fed into a recurrent neural network (RNN) or a Long Short-Term Memory (LSTM) network, which generates a satisfactory English description of the image.



ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.429 Volume 9 Issue VII July 2021- Available at www.ijraset.com

#### III. BACKGROUND AND PROPOSED APPROACH

In this section, we look at a deep learning model that generates image descriptions automatically. Our model implements a CNN to extract image features, which are then fed into an LSTM network to predict a valid English description of the image. We can also use the Web Speech API's Speech Synthesis Utterance interface to convert the LSTM model's descriptions into speech.



Figure 3: Proposed Model

#### A. CNN-based Image Feature Extractor

For feature extraction, we use a CNN. CNNs have been extensively studied and used for image tasks, and they are now the most advanced methods for object recognition and detection [20]. In particular, we extract features from all input images using MobileNet, a simple but efficient and computationally light convolutional neural network for mobile vision applications. Object detection, fine-grained classifications, localization and face attributes, are just a few of the real-world applications that use MobileNet. In this we will explain the overview of MobileNet and how exactly it becomes the most efficient and lightweight neural network. And due to the computational constraints, we also reduced a 4096-Dim image feature vector to 512-Dim image feature vector using Principal Component Analysis (PCA). We give these features as an input into our LSTM at the first iteration.[27]

An input layer, hidden layers, and an output layer make up a convolutional neural network. A feed-forward neural network's middle layers are referred to as hidden because the activation function and final convolution cover their inputs and outputs. The hidden layers of a convolutional neural network include convolutional layers. This usually consists of a multiplication or other scalar product layer with ReLU as its activation function. Other convolution layers, such as pooling layers, normalisation layers and fully connected layers, follow.



Figure 4: The layers of a convolutional neural network

The CNN has three most-used layers: convolution, pooling and fully connected layers. Also, Rectified Linear Units (ReLU) functions as given below are employed because of the non-linear active function. The ReLU is faster than the traditional tanh function.

$$f(x) = \text{ReLU}(x) = \max(0, x)$$
(1)  
$$f(x) = \tanh(x) = (1 + e^{-x}) - 1$$
(2)

Dropout layer is used to prevent overfitting. With a probability of 50%, the dropout sets the output of every hidden neuron to zero (i.e., 0.5). The "dropped out" neurons don't engage in back propagation and don't contribute to the pass.



## International Journal for Research in Applied Science & Engineering Technology (IJRASET) ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.429

Volume 9 Issue VII July 2021- Available at www.ijraset.com

In comparison to other image classification algorithms, CNNs require very little pre-processing. This means that the network picks up on the filters that were previously hand-crafted in traditional algorithms. In feature design, this lack of reliance on prior knowledge and human effort could be a significant benefit.

#### B. LSTM-based Sentence Generator

We will use LSTM (Long Short Term Memory), a type of RNN (recurrent neural network) that is compatible with sequence prediction problems, to predict what the next words will be based on the previous text. It has outperformed traditional RNNs by overcoming the limitations of RNNs with short term memory. With a forget gate, LSTM can perform relevant information throughout the processing of inputs while discarding non-relevant information.[27]

Our LSTM model takes the image I and a sequence of input vectors (i1, ..., iT). Then it computes a sequence of hidden states (hs1, ..., hst) and a sequence of outputs (o1, ..., ot) by following the recurrence relation for t = 1 to T:

$$\begin{split} bv &= W_{hi}[CNN(I)] \eqno(1) \label{eq:while} \\ hst &= f(W_{hx}i_t + W_{hh}hs_t - 1 + b_h + 1(t=1) \circ b_v) \eqno(2) \\ ot &= Softmax(W_{oh}hs_t + b_o) \eqno(3) \end{split}$$

where  $W_{hsi}$ ,  $W_{hsx}$ ,  $W_{hsss}$ ,  $W_{ohs}$ ,  $x_i$ ,  $b_{hs}$ , and  $b_o$  are learnable parameters and CNN(I) represents the image features or object extracted by the CNN model..



Figure 5: LSTM unit and its gates

The probability of each word in the vocabulary is the LSTM's output.

C. Speech Synthesis Utterance (For Converting Caption to Speech)

Speech synthesis is the process of creating a human voice. To turn the description into audio, we use the Speech Synthesis Utterance API. The Speech Synthesis Utterance is a Web Speech API interface that represents a speech request. It contains the content that the speech service should read as well as for instructions on how to read it (e.g., language, pitch, and volume.)[25].



Figure 6: Text to Speech



ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.429 Volume 9 Issue VII July 2021- Available at www.ijraset.com

#### IV. EXPERIMENTS

#### A. Dataset

There are a few datasets utilized for preparing, testing, and assessment of the picture subtitling techniques. The datasets contrast in different points of view like the quantity of pictures, the quantity of inscriptions per picture, the configuration of the subtitles, and picture size. Three datasets which are Flickr8k, Flickr30k, and MS COCO Dataset are prevalent.

For this paper, we will utilize the Flickr 8K dataset which has become the standard testbed for picture subtitling. There additionally are other enormous datasets like Flickr\_30K and MSCOCO dataset yet it can require weeks just to mentor the organization so we'll utilize a little Flickr 8k dataset. The benefit of a gigantic dataset is that we will construct better models. The dataset comprises 8,000 preparing pictures and their description. By partnering each picture with multiple, independently produced sentences, the dataset captures a number of the linguistic variety which will be wont to describe an equivalent image.

Flickr 8k is a decent beginning dataset as it is little and can be prepared effectively on low-end workstations/work areas.

| Dataset Filename         | Description   |
|--------------------------|---|
| Flick8k_Dataset          | It contains 8000 images                             |
| Flickr8k.token.txt       | It contains the image id along with the 5 captions. |
| Flickr8k.trainImages.txt | It contains the training image id's.                |
| Flickr8k.testImages.txt  | It contains the test image id's.                    |

 Table 1: Dataset Structure

#### B. Training

Based on the current word (xt) and the preceding context, we train our LSTM model to accurately predict the next word (ot) (ht1). This is how we go about it: We set hs0 to 0, i1 to the START vector, and o1 to the first word in the sequence as the desired label. The first word generated by the network is then assigned to the word vector i2. Based on the prior context, the network predicts the first word, the second word, and so on. iT represents the last word, and oT is set to an END token in the final step.



Figure 7: Visual representation of the final model.



#### C. Testing

To predict a sentence, we acquire the picture highlights by, set h0 to the 0, set x1 to the START vector, and figure the appropriation over the principal word y1. Appropriately, we pick the argmax from the circulation, set its inserting vector as x2, and rehash the system until the END token is produced.

| Layer (type)                               | Output | Shape    | Param # | Connected to  |
|--|--------|----------|---------|---|
| input_3 (InputLayer)                       | (None, | 40)      | 0       |   |
| <pre>input_2 (InputLayer)</pre>            | (None, | 1000)    | 0       |   |
| embedding_1 (Embedding)                    | (None, | 40, 64)  | 528384  | input_3[0][0]   |
| dense_1 (Dense)                            | (None, | 64)      | 64064   | input_2[0][0]   |
| lstm_1 (LSTM)                              | (None, | 40, 128) | 98816   | embedding_1[0][0]   |
| <pre>repeat_vector_1 (RepeatVector)</pre>  | (None, | 40, 64)  | 0       | dense_1[0][0]   |
| <pre>time_distributed_1 (TimeDistrib</pre> | (None, | 40, 64)  | 8256    | lstm_1[0][0]  |
| concatenate_1 (Concatenate)                | (None, | 80, 64)  | 0       | <pre>repeat_vector_1[0][0] time_distributed_1[0][0]</pre> |
| bidirectional_1 (Bidirectional)            | (None, | 256)     | 197632  | concatenate_1[0][0]                                       |
| dense_3 (Dense)                            | (None, | 8256)    | 2121792 | bidirectional_1[0][0]                                     |
| activation_1 (Activation)                  | (None, | 8256)    | 0       | dense_3[0][0]   |
| Total papamer 2 010 044                    |        |          |         |   |

Total params: 3,018,944 Trainable params: 3,018,944

Non-trainable params: 0

#### Figure 8: Model Summary

#### RESULTS

V.

Our models generate sensible descriptions of images in valid English (Figure 9) with audio of the descriptions. The model identifies visual-semantic correspondences that can be understood. The generated descriptions are accurate enough for blind or visually impaired people to use. In general, we find that the training data contains a significant portion of the generated sentences.



The man at bat readies to swing at the pitch while the umpire looks on.



A horse carrying a large load of hay and two people sitting on it.



A large bus sitting next to a very tall building.



Bunk bed with a narrow shelf sitting underneath it.



ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.429 Volume 9 Issue VII July 2021- Available at www.ijraset.com



A person is walking along a beach with a big dog



Ar .

black and white dog carries tennis ball in its mouth

A man is doing a trick on a snowboard



Figure 9: Example images and captions from the Flicker8k Caption dataset

A

а



Figure 10: Results of images that give description by generating accurate captions.



| Image Caption Generator x +   | - o ×       |
|---|-------------|
| $\leftrightarrow \rightarrow C$ (O localhost:800  | 📩 🖈 🔮 💊 💘 🔤 |
| INDERPRASTHA ENGINEERING COLLEGE<br>Approved by AICTE & Affiliated to Dr. APJ Abdul Kalam Technical University, Lucknow, U.P.)<br>College Code-0030   | T S<br>S T  |
| IMAGE CAPTION GENERATOR<br>This is a Assistive Vision Technology model based on CNN-LSTM neural networks which automatically detects the objects in the<br>images and generates descriptions for the images. These descriptions are converted to speech with the help Speech Synthesis<br>Utterance interface of the Web Speech API from which it can be read out loud to the person. | 3           |
| Good to go!   |             |
| Upload Your Picture Here  |             |
| Upload Images   |             |
| Generate Caption  |             |
| Caption:  |             |

Figure 11: UI of Web Page



Figure 12: Example 1



Figure 13: Example 2

#### VI. CONCLUSION

In practically every complicated area of AI, image captioning has several benefits. The main use case of our model is to assist the visually impaired to know the environment and make it easy to act consistently with the environment. As this is often a posh task to try to do, with the assistance of pre-trained models and powerful deep learning frameworks like TensorFlow and Keras, we made it possible. This is completely a Deep Learning project which makes use of multiple Neural Networks like Convolutional Neural networks and LSTM to detect objects and caption the pictures. To deploy our model as a web application, we have used HTML, CSS and JavaScript as frontend and NodeJS as the backend

We presented a deep learning-based model for automatically generating image descriptions with the goal of assisting visually impaired people in better understanding their surroundings. Our described model is based on a MobileNet which is CNN pre-trained model that encodes a picture into

a compact representation, followed by an LSTM that generates corresponding sentences supported by the learned image features. We showed that this model achieves state-of-the-art performance, in which the generated captions are highly descriptive of the entities and scenes depicted in the pictures. Using text-to-speech technology, visually impaired people can considerably benefit and have a far better sense of their surroundings due to the excellent quality of the generated image descriptions. Our current approach only creates captions for images, which is a difficult work in and of itself, and captioning live video frames is even more challenging. It is often completely GPU-based and captioning live video frames is not possible with the overall CPUs. Video captioning is a popular study subject that is changing people's lives, with application cases found in practically every domain. In the following years, we hope to accomplish this aim of video captioning to broaden the scope.



ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.429 Volume 9 Issue VII July 2021- Available at www.ijraset.com

#### REFERENCES

- [1] "Blindness and Visual Impairment." WHO stands for the World Health Organization (2014). 10 April 2016. Web.
- [2] Jia Deng, Hao Su, Sanjeev Satheesh, Sean Ma, Olga Russakovsky, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Jonathan Krause, Alexander C. Berg, and Li Fei-Fei. "ImageNet Large Scale Visual Recognition Challenge." Russakovsky, Olga, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Sa Web. 19 April 2016. Int J Comput Vis 115.3 (2015): 211-52.
- [3] Andrew Zisserman, Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. Everingham, Mark, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. "In Pascal, the Visual Object Classes (VOC) Challenge." 303-38 in Int J Comput Vis 88.2 (2009). 22 May 2016. Web. The International Journal of Computer Vision is a publication dedicated to the study of computer vision 303-38 in Int J Comput Vis 88.2 (2009).
- [4] Ali, Mohsen Hejrati, Mohammad Amin Sadeghi, Peter Young, Cyrus Rashtchian, Julia Hockenmaier, and David Forsyth. Farhadi, Ali, Mohsen Hejrati, Mohammad Amin Sadeghi, Peter Young, Cyrus Rashtchian, Julia Hockenmaier, and David Forsyth. "Generating Sentences from Images: Every Picture Tells a Story." 15-29. Web. 5 April 2016. Computer Vision ECCV 2010 Lecture Notes in Computer Science (2010): 15-29.
- [5] Sagnik Dhar, Visruth Premraj, Girish Kulkarni, Siming Li, Yejin Choi, Alexander C. Berg, and Tamara L. Berg. Kulkarni, Girish, Visruth Premraj, Sagnik Dhar, Siming Li, Yejin Choi, Alexander C. Berg, and Tamara L. Berg. "Baby Talk: Understanding and Generating Simple Image Descriptions," as the title suggests. 2011 Cvpr (2011). 27 May 2016. Web.
- [6] Li, Li-Jia, R. Socher, and Li Fei-Fei. "Towards Total Scene Understanding: Classification, Annotation and Segmentation in an Automatic Framework." 2009 IEEE Conference on Computer Vision and Pattern Recognition (2009). Web. 21 Apr. 2016
- [7] Stephen Gould, Richard Fulton, and Daphne Koller are the authors of this book. "Geometrically and Semantically Consistent Regions in a Scene." IEEE's 12th International Conference on Computer Vision was held in 2009. (2009). 6 May 2016. Web.
- [8] Sanja Fidler, Abhishek Sharma, and Raquel Urtasun Fidler, Sanja, Abhishek Sharma, and Raquel Urtasun Fidler, Sanja, Abhishek "A Thousand Pixels Is Worth a Sentence." The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) was held in 2013. (2013). 18 May 2016 (web).
- [9] "What, Where, and Who?" Li Fei-Fei. Scene and Li Li-Jia and Object Recognition for Event Classification." The 11th IEEE International Conference on Computer Vision was held in 2007. (2007). 10 April 2016. Web.
- [10] Russakovsky et al.2015, Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. IJCV, 115(3):211–252.
- [11] David Eigen, Xiang Zhang, Michael Mathieu, Sermanet et al.2013 Pierre Sermanet, Rob Fergus, and Yann "LeCun. 2013. Overfeat: Integrated recognition, localization and detection using convolutional networks. arXiv preprint arXiv:1312.6229
- [12] Prashanth Mannem and Ankush Gupta From image annotation to image description, we've got you covered. In the processing of neural information.
- [13] Lisa Anne Hendricks, Jeffrey Donahue, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell are among those who have appeared on the show. Visual recognition and description using long-term recurrent convolutional networks. IEEE CVPR is a journal that publishes research papers.
- [14] Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko are Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko, respectively. IEEE ICCV video to text-sequence to sequence-sequence to computer Vision).
- [15] Dumitru Erhan, Oriol Vinyals, Alexander Toshev, Samy Bengio A neural image caption generator to show and tell. IEEE CVPR is a journal that publishes research papers.
- [16] Kuznetsova, Polina, Vicente Ordonez, Tamara Berg, Yejin Choi. "TREETALK: Composition and Compression of Trees for Image Descriptions." Transactions of the Association for Computational Linguistics 2 (2014): 351-362. Web. 1 Apr. 2016
- [17] Gupta and Mannem. "From image annotation to image description. In Neural information processing." Springer (2012). Web. 7 Apr. 2015
- [18] "Gradient-based learning applied to document recognition," by LeCun, Bottou, Bengio, and Haffner. IEEE Transactions on Industrial Electronics, vol. 86, no. 11, pp. 22782324, 1998. 27 May 2016. Web.
- [19] "Imagenet classification with deep convolutional neural networks," by Krizhevsky, Sutskever, and Hinton. NIPS is a non-profit organisation dedicated to (2012). The 28th of April, 2016 is the date on which this article was published on the internet
- [20] "Deep Visual-semantic Alignments for Generating Image Descriptions," by Andrej Karpathy and Li Fei-Fei. The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2015 is a gathering of experts in the field of computer vision and pattern recognition (2015). 29 May 2016. Web.
- [21] Ruslan Salakhutdinov, Kiros Ryan, and Rich Zemel 595-603 in Proceedings of the 31st International Conference on Machine Learning (ICML-14) (2014). Web. May 21st, 2016
- [22] Christopher Elamri, Teun de Planque Department of Computer Science Stanford University {mcelamri, teun}@stanford.edu
- [23] Researchers and educators who want to use the images for non-commercial research and educational purposes can use the Flickr 8k dataset. illinois.edu/sec/1713398
- [24] Text-to-Speech, developed by S. Venkateswarlu, Duvvuri B K Kamesh Duvvuri, and Sastry Jammalamadaka in 2016, can convert text images into sound with a good enough performance and a readability tolerance of less than 2%.
- [25] Srudeep PA gives an overview of MobileNet in 2020. A Mobile Vision CNN That Works.
- [26] "Long Short-Term Memory," by Sepp Hochreiter and Jrgen Schmidhuber. Web. 23 April 2016. Neural Computation 9.8 (1997): 1735-780.











45.98



IMPACT FACTOR: 7.129







# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 🕓 (24\*7 Support on Whatsapp)