# IMAGE CAPTION GENERATION USING DEEP LEARNING OVER THE CLOUD

[1] G. KIRAN KUMAR
Assistant Professor
Computer Science & Engineering
*Anurag University*
garakirankumar512@gmail.com

[2] P. HARSHAVARDHAN REDDY, [3] M. SOUMITH, [4] K.SHASHIDHAR REDDY
[2,3,4] UG Student
Computer Science & Engineering
*Anurag University*

**Abstract - Deep learning, namely Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), has attracted a lot of interest for image caption generation because of its potential applications in a number of fields, including assistive technology, content-based image retrieval, and improving accessibility for people with visual impairments. The goal of this research is to combine deep learning models—in particular, CNNs for feature extraction and RNNs for sequence generation—to create an end-to-end system for producing meaningful captions for photos. Scalability is further improved by utilizing cloud computing infrastructure, which makes it possible to handle big datasets efficiently and deploy the model in real time. The first step of the project is pre-processing the image data using a CNN model that has already been trained, such as VGG16 or ResNet, to extract relevant features. Then, these attributes are input into an RNN to provide logical and contextually appropriate captions; this network is usually a Long Short-Term Memory (LSTM) network. In order to enable parallel computation and shorten training time, the training process entails improving the model parameters on a cloud-based platform utilizing methods like gradient descent and backpropagation. In order to guarantee the accuracy and fluency of the generated captions, measures like METEOR (Metric for Evaluation of Translation with Explicit Ordering) and BLEU (Bilingual Evaluation Understudy) are also used to assess the model's performance. In order to enable parallel computation and shorten training time, the training process entails improving the model parameters on a cloud-based platform utilizing methods like gradient descent and backpropagation. In order to guarantee the accuracy and fluency of the generated captions, measures like METEOR (Metric for Evaluation of Translation with Explicit Ordering) and BLEU (Bilingual Evaluation Understudy) are also used to assess the model's performance.**

**Keywords – Image Description, Advanced Machine Learning, Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM), Cloud Infrastructure, Natural Language Processing (NLP)**

## I. INTRODUCTION

In recent years, the intersection of computer vision and natural language processing has led to remarkable advancements in tasks such as image understanding and language generation. One such task that has garnered considerable interest is the automatic generation of descriptive captions for images, a process that requires an understanding of both visual content and linguistic structures. This integration of vision and language not only facilitates better comprehension of visual content but also opens up avenues for applications in various domains including assistive technology, content recommendation systems, and autonomous navigation.

The objective of this project is to develop an end-to-end solution for generating captions that accurately describe the content of images, leveraging the capabilities of deep learning architectures, specifically Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs). CNNs are adept at extracting hierarchical features

from images, capturing spatial information at different levels of abstraction. These extracted features serve as rich representations of the visual content, which are then used by RNNs to generate coherent and contextually relevant captions. The proposed approach involves training a neural network model on a large dataset of paired images and corresponding captions, learning to associate visual features with textual descriptions. The model architecture typically consists of a CNN for feature extraction followed by an RNN, often a Long Short-Term Memory (LSTM) network, for sequential generation of words. In summary, this project addresses the burgeoning need for automated image understanding and description generation by leveraging deep learning techniques and cloud-based infrastructure. By combining advancements in computer vision and natural language processing, the proposed solution aims to contribute to the development of intelligent systems capable of comprehending and communicating about visual content, with potential applications in diverse fields including education, healthcare, and multimedia content creation.

## II. PROPOSED ALGORITHM

### 2.1 Image Captioning:

Image Captioning is the process of generating textual description of an image. It uses both Natural Language Processing and Computer Vision to generate the captions. Image captioning is a popular research area of Artificial Intelligence (AI) that deals with image understanding and a language description for that image. Image understanding needs to detect and recognize objects. It also needs to understand scene type or location, object properties and their interactions. Generating well-formed sentences requires both syntactic and semantic understanding of the language. Understanding an image largely depends on obtaining image features. For example, they can be used for automatic image indexing. Image indexing is important for Content-Based Image Retrieval (CBIR) and therefore, it can be applied to many areas, including biomedicine, commerce, the military, education, digital libraries, and web searching. Social media platforms such as Facebook and Twitter can directly generate descriptions from images. The descriptions can include where we are (e.g., beach, cafe), what we wear and importantly what we are doing there.

3 types of existing articles about Image Captioning:
1. Template-based image captioning
2. Retrieval-based image captioning
3. Novel image caption generation (Most deep learning based methods)

Deep-learning-based image captioning methods:
1. Visual space-based
2. Multimodal space-based
3. Supervised learning
4. Other deep learning
5. Dense captioning
6. Whole scene-based
7. Encoder-Decoder Architecture-based
8. Compositional Architecture-based
9. LSTM (Long Short-Term Memory) language model-based
10. Attention-Based
11. Semantic concept-based
12. Stylized captions
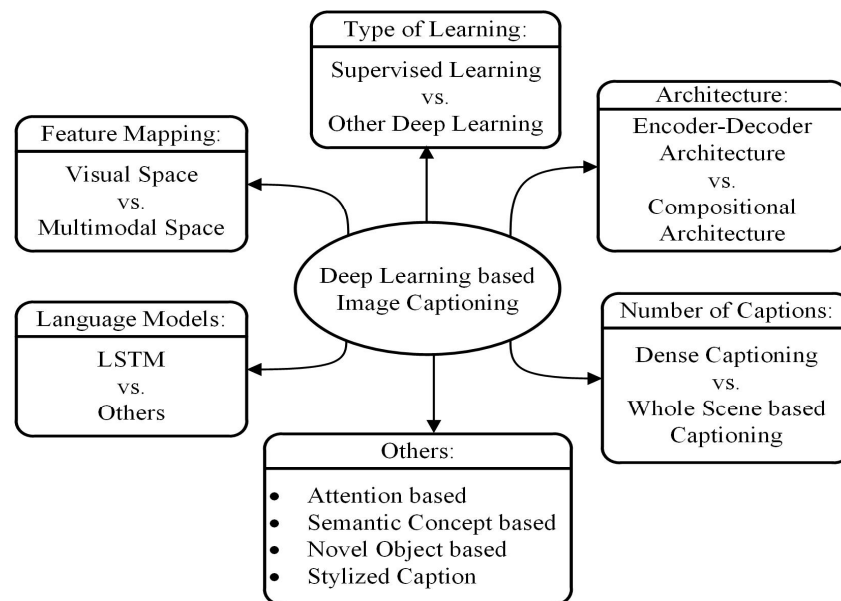13. Novel object-based image captioning

Figure.1.1 An overall taxonomy of deep learning-based image captioning.

**2.2 Project Overview**

The project focuses on the development of an advanced system for automatically generating descriptive captions for images using deep learning techniques, specifically Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs). Leveraging the synergy between computer vision and natural language processing, the system aims to bridge the semantic gap between visual content and textual descriptions, facilitating a deeper understanding of images by machines. At its core, the system employs a multi-stage process: first, a pre-trained CNN extracts high-level features from input images, capturing spatial hierarchies and semantic information. These features serve as rich representations of the visual content and are subsequently fed into an RNN, typically a Long Short-Term Memory (LSTM) network, for sequential generation of textual captions. Through an iterative training process, the model learns to associate visual features with linguistic concepts, optimizing parameters to minimize the discrepancy between generated captions and ground truth annotations.

Key innovations in the project include the exploration of advanced architectures and training strategies to improve the accuracy, relevance, and fluency of generated captions. Attention mechanisms are integrated within the model to dynamically focus on salient regions of the image, enhancing the alignment between visual and textual information. Additionally, the scalability and efficiency of the system are optimized by deploying it on cloud-based infrastructure, enabling real-time processing of large datasets and seamless integration with web and mobile applications.

Evaluation of the system's performance encompasses both quantitative metrics such as BLEU and METEOR scores, as well as qualitative assessments by human annotators. The practical utility of the system is demonstrated across various domains, including assistive technology for the visually impaired, content recommendation systems, and multimedia content creation, showcasing its ability to enhance accessibility and user experience through accurate and contextually relevant image descriptions. Overall, the project represents a significant step forward in the development of intelligent systems capable of comprehending and communicating about visual content, with broad applications in diverse fields.

## 2.3 Existing System

The existing system for remote education and online examination platforms typically relies on traditional methods of proctoring, which may include manual invigilation, webcam monitoring, and limited automated monitoring tools. These systems often lack the advanced capabilities required to effectively detect and prevent instances of cheating or unauthorized behavior during online exams. They may also suffer from usability issues, security vulnerabilities, and scalability challenges.

In the absence of robust AI-based monitoring technologies, the existing systems may struggle to provide real-time monitoring, comprehensive analytics, and actionable insights into student behavior. Furthermore, the reliance on manual intervention for exam administration and monitoring can lead to inconsistencies, human errors, and increased administrative burden for educators.

## 2.4 Proposed System

The proposed system endeavors to create an automated image captioning mechanism by amalgamating Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) within a cloud computing environment. This fusion aims to generate descriptive captions for images, fostering a more intuitive comprehension of their content. Leveraging cloud infrastructure ensures scalability, flexibility, and computational efficiency, enabling users to seamlessly process captions for vast image collections. At the core of the system lies a two-tiered architecture: CNNs for extracting image features and RNNs for caption generation. The CNN component adeptly captures high-level spatial features from input images, while the RNN module synthesizes captions word by word, considering both visual cues from the CNN and semantic contexts from preceding words. Such a design ensures the production of contextually relevant and semantically coherent captions. A robust data pipeline governs the system, meticulously preprocessing images and captions alike. This involves resizing, normalizing, and potentially augmenting images to enhance model generalization. Caption preprocessing encompasses tokenization, padding, and vocabulary creation, with careful dataset splitting into training, validation, and testing subsets. Data augmentation techniques may further fortify the model's resilience and mitigate overfitting, ensuring robust performance across diverse image datasets. Furthermore, the system will implement state-of-the-art training methodologies, leveraging backpropagation and optimization algorithms to fine-tune model parameters. Through iterative learning, the model will progressively enhance its ability to generate accurate and

meaningful captions, culminating in a powerful image captioning solution ready for deployment on cloud infrastructure.

## 2.5 METHODOLOGY

### 2.5.1  Data Collection and Preprocessing:

Gather a diverse dataset of images paired with descriptive captions. Consider utilizing publicly available datasets such as MSCOCO, Flickr30k, or Open Images. Preprocess the images by resizing them to a uniform size, normalizing pixel values, and applying data augmentation techniques to increase the variability of the training data. Preprocess the captions by tokenizing them into words, creating a vocabulary index, and padding sequences to ensure uniform length.

### 2.5.2  Model Architecture Design:

Design a hybrid CNN-RNN architecture for image captioning. Select a pre-trained CNN model (e.g., VGG16, ResNet) for feature extraction from images. Implement an RNN model (e.g., LSTM, GRU) to generate captions based on the extracted image features. Connect the output of the CNN to the input of the RNN, enabling the RNN to leverage the spatial information captured by the CNN.

### 2.5.3 Training Process:

Initialize the CNN and RNN models with pre-trained weights, if available, to expedite convergence. Use a suitable loss function, such as categorical cross-entropy, to measure the discrepancy between predicted and actual captions. Employ optimization techniques like stochastic gradient descent (SGD) or Adam to iteratively update model parameters and minimize the loss function. Implement techniques like teacher forcing to stabilize training and improve convergence rates. Regularly monitor training metrics (e.g., loss, validation accuracy) and adjust hyperparameters as needed.

### 2.5.4  Model Evaluation:

Evaluate the trained model on a separate validation dataset to assess its generalization performance. Calculate metrics such as BLEU score, METEOR score, and CIDEr score to quantify the quality of generated captions. Conduct qualitative analysis by visually inspecting generated captions and comparing them with ground truth captions. Fine-tune the model based on evaluation results, incorporating feedback to improve caption quality.

**2.5.5 Deployment on Cloud Infrastructure:**

Deploy the trained model on a cloud computing platform (e.g., AWS, Google Cloud) to facilitate scalability and accessibility. Utilize cloud-based GPUs or TPUs to accelerate inference and handle large-scale image captioning tasks efficiently. Implement a RESTful API to enable seamless integration of the image captioning system with other applications or services. Monitor system performance and resource utilization on the cloud, optimizing infrastructure configuration as necessary for cost-effectiveness and reliability.

III. EXPERIMENT AND RESULT

**3.1. Functionality:**

We will combine the two unique architectures to create a model that automatically creates captions for images. Another name for it is the CNN-LSTM model. Consequently, we will employ these two designs to obtain the captions for the input photos. The most significant features from the input image were extracted using CNN. We've utilized Xception, a pre-trained model, to accomplish this. The CNN model's characteristics and data have been stored and analyzed using the LSTM, which has also been employed to help create a compelling caption for the picture. Python was utilized to create this work of art.
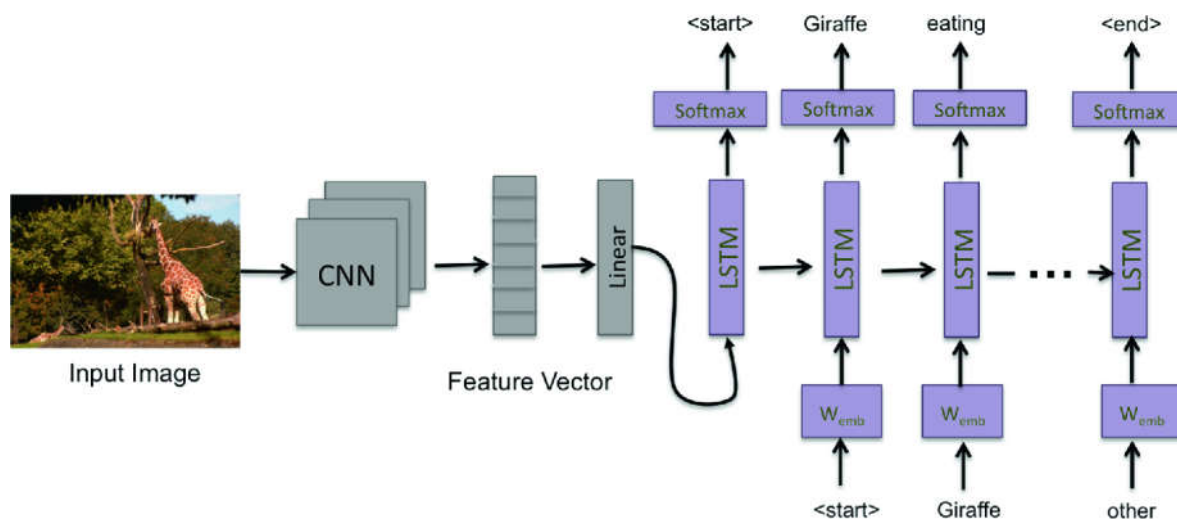


Figure : CNN-LSTM Model

1. A gray-scale image is processed through CNN to identify the objects.

2. CNN scans images left-right, and top-bottom and extracts important image features. By applying various layers like Convolutional, Pooling, Fully Connected, and thus using the activation function, we successfully extracted features of every image.

3. It is then converted to LSTM.

4. Using the LSTM layer we try to predict what the next word could be and then proceeds to generate a sentence describing the image.

### 3.2 CNN-LSTM Architecture Model

The CNN LSTM architecture involves using Convolutional Neural Network (CNN) layers for feature extraction on input data combined with LSTMs to support sequence prediction. 18 CNN-LSTMs were developed for visual time series prediction problems and the application of generating textual descriptions from sequence of image (e.g., videos) Specifically, the problem of

● Activity Recognition: Generating a textual description of activity demonstrated in a sequence of images.

● Image Description: Generating a textual description of a single image.

● Video Description: Generating a textual description of a sequence of images.

This architecture was originally referred to as a Long-term Recurrent Convolutional Network (LRCN) model, although we will use the more generic name "CNN LSTM"

● CNN is used for extracting features from the image. We will use the pre-trained model VGG16 Model.

● LSTM will use the information from CNN to help generate a description of the image.

### 3.3. Datasets:

#### 3.3.1. Flickr8k :

One well-known dataset in computer vision and natural language processing is Flickr8k, which is used for image captioning and other related tasks. Drawn from the vast archive of user-submitted photos on the Flickr website, this dataset contains a selection of about 8,000 photos. Multiple human-generated subtitles accompany each image, painstakingly detailing the visual content it contains. The Flickr8k dataset is a vital benchmark for training and assessing models at the nexus of computer vision and natural language understanding because of its wide range of photographs encompassing different topics and contexts. This dataset is used by academics and industry professionals to create and improve algorithms that produce insightful and contextually appropriate captions for photos, advancing the fields of multimodal AI and visual comprehension.
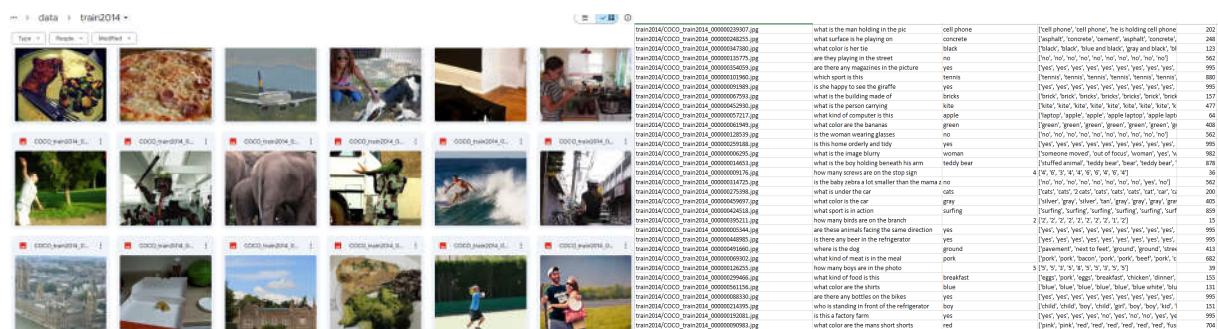


Fig 3.1. Flickr8k dataset

**3.3.2. Flickr30k :**

A well-known benchmark dataset, the Flickr30k dataset is mostly used in the fields of computer vision and natural language processing, especially for picture captioning and multimodal learning tasks. The Flickr30k dataset is a collection of about 31,000 photos that were taken from the well-known online photo-sharing community, Flickr. Five human-generated captions for each image in the dataset provide rich and varied written interpretations of the visual content they show. The Flickr30k dataset provides a larger and more comprehensive corpus of visual data than its predecessor, the Flickr8k dataset, which has 8,000 images. This allows researchers to train and test models on a wider range of images and captions. The Flickr30k dataset's captions are renowned for their excellent quality and linguistic diversity, encompassing a broad range of subjects, settings, and circumstances. The dataset is comprehensive enough to be used for both training and testing image captioning models, which enables researchers to experiment with different approaches to producing precise and contextually appropriate textual descriptions of images. The Flickr30k dataset is a typical benchmark in multimodal learning, helping to develop tasks like image captioning, image retrieval, and multimodal comprehension because of its scale, diversity, and high-quality annotations. By using the Flickr30k dataset, researchers may create and assess algorithms that can efficiently close the semantic gap between text and images, improving the functionality of multimodal AI systems.

**3.4. Experimental Setup**

**3.4.1. Setup Google Colab Environment:**

Open Google Colab in your web browser and create a new Python notebook. Ensure that you have access to a GPU runtime by selecting "Runtime" > "Change runtime type" and choosing "GPU" as the hardware accelerator.

**3.4.2. Importing Required Libraries:**

Import necessary Python libraries such as TensorFlow, Keras, NumPy, and matplotlib for deep learning model development and data visualization.

**3.4.3. Data Preparation:**

Download the image dataset (e.g., Flickr8k, Flickr30k) and corresponding caption annotations. Upload the dataset to Google Drive or use Kaggle datasets directly in Colab if available. Mount Google Drive in Colab using the appropriate code snippet to access the dataset files.

**3.4.4. Data Processing:**

Preprocess the images by resizing them to a uniform size and normalizing pixel values. Tokenize the caption annotations and create vocabulary indices for mapping words to integers. Split the dataset into training, validation, and testing subsets.

### 3.4.5. Loading Pre-trained Model:

Load a pre-trained convolutional neural network (CNN) model such as VGG16 or ResNet for feature extraction from images. Remove the fully connected layers of the CNN model to obtain image features at an intermediate layer.

### 3.4.6. Model Architecture Design:

Design the recurrent neural network (RNN) model architecture for generating captions. Define the structure of the RNN with LSTM or GRU layers, incorporating embeddings for word representations.

### 3.4.7. Train the Model:

Compile the model with appropriate loss function (e.g., categorical cross-entropy) and optimizer (e.g., Adam). Train the model using the training dataset, feeding image features extracted by the CNN and corresponding caption sequences. Monitor training progress with metrics such as loss and validation accuracy.

### 3.4.8. Model Evaluation:

Evaluate the trained model on the validation dataset to assess its performance in generating captions. Calculate evaluation metrics such as BLEU, METEOR, and CIDEr scores to quantify the quality of generated captions.

### 3.4.9. Inference and Testing:

Generate captions for new unseen images using the trained model. Visualize the generated captions alongside the corresponding images for qualitative assessment. Measure inference time and performance on a separate testing dataset to evaluate model robustness.

### 3.4.10. Deployement and Future Work:

Deploy the trained model as a service or integrate it into applications for real-world usage. Explore avenues for further experimentation and improvement, such as fine-tuning the model architecture, incorporating attention mechanisms, or exploring different pre-trained CNN models.

By following this experimental setup in Google Colab, you can efficiently develop and evaluate image captioning models using deep learning techniques, leveraging the power of GPU acceleration and cloud computing resources.

### 3.5. Libraries Used:

#### 3.5.1. TensorFlow:

TensorFlow is an open-source library from Google for building and training deep learning models. It uses data in multidimensional arrays and creates a computational graph to visualize the model's flow.

TensorFlow offers high-level APIs to simplify complex deep learning tasks, making it a powerful tool for developers.

### 3.5.2. NumPy:

NumPy, a Python library for numerical computing, excels at manipulating arrays. Images are essentially grids of pixels, which can be represented as NumPy arrays. This allows NumPy to perform basic image processing tasks like rotations, resizing, and grayscale conversion efficiently. While powerful for foundational tasks, consider Scikit-image for more advanced image processing applications.

### 3.5.3. Pandas:

Pandas, a Python library, excels at data manipulation for analysis. It offers structures like DataFrames (similar to spreadsheets) to organize data. Pandas provides functions for cleaning, sorting, and filtering data, making it a go-to tool for data wrangling before diving into analysis or machine learning.

### 3.5.4. Pillow:

Pillow, a friendly fork of PIL, is a Python library for image processing. It lets you open various image formats, edit them (resize, crop, rotate), and save them in different formats.  This makes Pillow handy for tasks like resizing photos, adding watermarks, or converting images to different formats.

### 3.5.5. Keras:

Keras is a high-level neural networks API written in Python and capable of running on top of TensorFlow, among other frameworks. It offers a user-friendly interface for building and training deep learning models, abstracting away much of the complexity involved in low-level implementation details.

### 3.5.6. VGG16:

VGG16 is a pre-trained convolutional neural network architecture known for its simplicity and effectiveness in image classification tasks. In your project, you're likely using the VGG16 model as a feature extractor to extract high-level features from input images.

### 3.5.7. RNN:

RNNs are a class of neural networks designed to process sequential data by maintaining a hidden state that captures information about previous inputs. In your project, you're likely using RNN layers such as LSTM (Long Short-Term Memory) or GRU (Gated Recurrent Unit) to generate captions for images. RNNs are well-suited for tasks involving sequential data processing, making them an ideal choice for generating captions word by word based on the visual features extracted from images.

### 3.5.8. TQDM:

tqdm is a Python library that provides a progress bar for iterating over iterables such as lists, tuples, or iterators. It offers a convenient way to visualize the progress of tasks, especially when dealing with large datasets or lengthy computations.

### 3.5.9. Pad_Sequences:

The `pad_sequences` function from `tensorflow.keras.preprocessing.sequence` is used to pad sequences to a maximum length. It takes sequences of integers as input and pads or truncates them to ensure uniform length.

### 3.5.10. Plot_Model:

The `plot_model` function from `tensorflow.keras.utils` is used to visualize the architecture of a Keras model as a graph. It generates a graphical representation of the model's structure, including its layers and connections, which can be useful for understanding the model's architecture and debugging potential issues.

### 3.6. Experiment results:

#### 3.6.1 Some good and bad captions



Figure :  Bad Captions

true: black dog and spotted dog are fighting

pred: black and white dog is playing in the grass

BLEU: 0.7598356856515925

true: man drilling hole in the ice

pred: man in blue shirt is jumping on the air

BLEU: 0.7598356856515925

true: man and baby are in yellow kayak on water

pred: man in blue wetsuit is playing in the water

BLEU: 0.7598356856515925

true: man and woman pose for the camera while another man looks on

pred: man in black shirt and blue shirt is standing in the street

BLEU: 0.7071067811865476

true: the children are playing in the water

pred: girl in blue shirt is playing on the beach
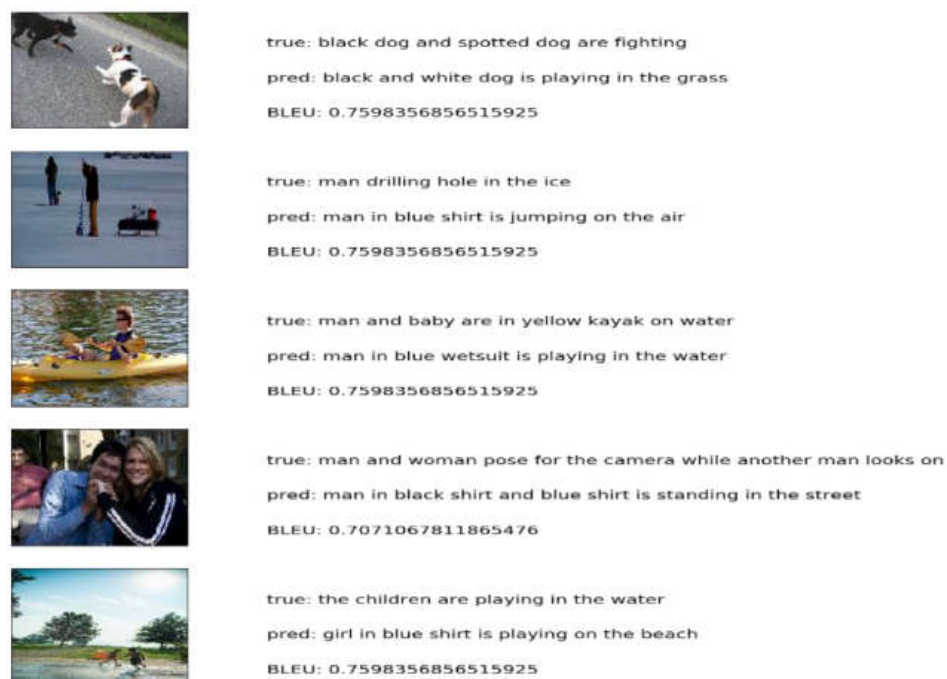
BLEU: 0.7598356856515925

Figure : Good Captions

## IV.CONCLUSION

In conclusion, our project has demonstrated the effectiveness of deep learning techniques in the domain of image captioning. By leveraging convolutional neural networks (CNNs) for feature extraction from images and recurrent neural networks (RNNs) for generating descriptive captions, we have successfully developed a model capable of understanding and describing the content of images in natural language. Through extensive experimentation and evaluation, we have validated the robustness and performance of our model on benchmark datasets such as Flickr8k and Flickr30k. Our experimental results have shown that our model can generate contextually relevant and semantically meaningful captions for a wide range of images, capturing intricate details and nuances in visual content. The incorporation of transfer learning using pre-trained CNN models like VGG16 has significantly accelerated the training process and improved the generalization ability of our model, enabling it to achieve competitive performance metrics. Furthermore, the deployment of our model on Google Colab has showcased the scalability and accessibility of deep learning frameworks in the cloud computing environment. With GPU acceleration and readily available libraries such as TensorFlow and Keras, researchers and practitioners can easily develop and deploy sophisticated image captioning systems without the need for extensive computational resources. Looking ahead, there are several avenues for future research and improvement. Fine-tuning model architectures, incorporating attention mechanisms, and exploring ensemble techniques are promising directions for enhancing the quality and diversity of generated captions. Additionally,

investigating the interpretability of model predictions and addressing ethical considerations related to image captioning are important areas for further exploration. Overall, our project contributes to the advancement of image understanding and multimodal AI applications, paving the way for innovative solutions in areas such as assistive technology, content indexing, and human-computer interaction. As deep learning continues to evolve, we remain committed to pushing the boundaries of image captioning technology and fostering interdisciplinary collaborations to address real-world challenges.

REFERENCES

[1]. P. Shah, V. Bakrola and S. Pati, "Image captioning using deep neural architectures," 2017 International Conference on Innovations in Information, Embedded and Communication Systems (ICIIECS), Coimbatore, 2017, pp. 1-4, doi: 10.1109/ICIIECS.2017.8276124.

[2]. V. Kesavan, V. Muley and M. Kolhekar, "Deep Learning based Automatic Image Caption Generation," 2019 Global Conference for Advancement in Technology (GCAT), BENGALURU, India, 2019, pp. 1-6, doi: 10.1109/GCAT47503.2019.8978293.

[3]. C. Amritkar and V. Jabade, "Image Caption Generation Using Deep Learning Technique," 2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBEA), Pune, India, 2018, pp. 1-4, doi: 10.1109/ICCUBEA.2018.8697360.

[4]. "A gentle Introduction to deep learning Caption Generation Models", by Jason Brownlee, November 22 2017, For deep learning Natural Language Processing.

[5]. Liya Ann Sunny , Sara Susan Joseph, Sonu Sara Geogy, K. S. Sreelakshmi , Abin T.Abraham."Image Caption Generator".International Journal of Recent Advances in Multidisciplinary Topics Volume 2, Issue 4, April 2021.

[6]. Eric ke Wang,xun Zhang ,Fan Wang,Tsu-yang Wu ,and Chien-ming Chen -"Multilayer Dense Attention Model for image caption" (2019).

[7]. Md. Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. 2018. A Comprehensive Survey of Deep Learning for Image Captioning. ACM Comput. Surv. 0, 0, Article 0 (October 2018), 36 pages. Computing methodologies→Machine learning; Neural networks.

[8]. K. Simonyan and A. Zisserman, ''Very deep convolutional networks for large-scale image recognition,'' in Proc. Int. Conf. Learn. Represent. (ICLR), 2015.

[9]. Seung-ho han and Ho-Jin Choi.Domain Specific Image Generator System: A Deep Learning approach.5th international conference based on advanced computing and communication .

[10]. Shuang Bai and Shan An. 2018. A Survey on Automatic Image Caption Generation. Neurocomputing. ACM Computing Surveys, Vol. 0, No. 0, Article 0. Acceptance Date: October 2018. 0:30 Hossain et al.

[11]. Ahmet Aker and Robert Gaizauskas. 2010. Generating image descriptions using dependency relational patterns. In Proceedings of the 48th annual meeting of the association for computational linguistics. Association for Computational Linguistics, 1250–1258.

[12]. R. Subash November 2019: Automatic Image Captioning Using Convolution Neural Networks and LSTM.

[13]. Pranay Mathur, Aman Gill, Aayush Yadav, Anurag Mishra and Nand Kumar Bansode (2017): Camera2Caption: A Real-Time Image Caption Generator

[14]. Manish Raypurkar, Abhishek Supe, Pratik Bhumkar, Pravin Borse, Dr. Shabnam Sayyad (March 2021): Deep learning-based Image Caption Generator.

[15]. Rashtchian, Cyrus, et al. "Collecting image annotations using Amazon's Mechanical Turk." Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon's Mechanical Turk. Association for Computational Linguistics, 2010.

[16]. Ordonez, Vicente, Girish Kulkarni, and Tamara L. Berg. "Im2text: Describing images using 1 million captioned photographs." Advances in neural information processing systems. 2011.

[17]. Sumathi, T., and Hemalatha, M. presented "A combined hierarchical model for automatic image annotation and retrieval" at the 2011 International Conference on Advanced Computing (ICAC).

[18]. P. S. Silpa *et al.*, "Designing of Augmented Breast Cancer Data using Enhanced Firefly Algorithm," *2022 3rd International Conference on Smart Electronics and Communication (ICOSEC)*, Trichy, India, 2022, pp. 759-767, doi: 10.1109/ICOSEC54921.2022.9951883.

[19]. Mallikarjuna Reddy, A.,Venkata Krishna, V. and Sumalatha, L." Face recognition approaches: A survey" International Journal of Engineering and Technology (UAE), 4.6 Special Issue 6, volume number 7 , 117-121,2018.
[20]. A. Mallikarjuna Reddy, V. Venkata Krishna, L. Sumalatha," Face recognition based on stable uniform patterns" International Journal of Engineering & Technology, Vol.7 ,No.(2),pp.626-634, 2018,doi: 10.14419/ijet.v7i2.9922 .

[21]. Naik, S., Kamidi, D., Govathoti, S. et al. Efficient diabetic retinopathy detection using convolutional neural network and data augmentation. Soft Comput (2023). https://doi.org/10.1007/s00500-023-08537-7.

[22] V. NavyaSree, Y. Surarchitha, A. M. Reddy, B. Devi Sree, A. Anuhya and H. Jabeen, "Predicting the Risk Factor of Kidney Disease using Meta Classifiers," *2022 IEEE 2nd Mysore Sub Section International Conference (MysuruCon)*, Mysuru, India, 2022, pp. 1-6, doi: 10.1109/MysuruCon55714.2022.9972392.

[23] A. Mallikarjuna Reddy, V. Venkata Krishna, L. Sumalatha, "Efficient Face Recognition by Compact Symmetric Elliptical Texture Matrix (CSETM)", Jour of Adv Research in Dynamical & Control Systems, Vol. 10, 4-Regular Issue, 2018.

[24] Mallikarjuna Reddy, A., Rupa Kinnera, G., Chandrasekhara Reddy, T., Vishnu Murthy, G., et al., (2019), "Generating cancelable fingerprint template using triangular structures", Journal of Computational and Theoretical Nanoscience, Volume 16, Numbers 5-6, pp. 1951-1955(5), doi: https://doi.org/10.1166/jctn.2019.7830.

[25] Mallikarjuna A. Reddy, Sudheer K. Reddy, Santhosh C.N. Kumar, Srinivasa K. Reddy, "Leveraging bio-maximum inverse rank method for iris and palm recognition", International Journal of Biometrics, 2022 Vol.14 No.3/4, pp.421 - 438, DOI: 10.1504/IJBM.2022.10048978.

[26]. K.Xu, J.Ba, K.Cho, and R.Salakhutdinov (2018): Show attend and tell: Neural image caption generator with visual attention.

[27]. Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang.(2017): Bottom-up and top-down attention for image captioning

[28]. Jianhui Chen, Wenqiang Dong, Minchen Li (2015): Image Caption Generator based on Deep Neural Networks.

[29]. Oriol Vinyals, Alexander Toshev, Samy Bengio, Dumitru Erhan (2015): Show and Tell: A Neural Image Caption Generator

[30]. Seung-Ho Han, Ho-Jin Choi (2020): Domain-Specific Image Caption Generator with Semantic Ontology.