

IMAGE CAPTION GENERATOR USING DEEP LEARNING

¹ SAHANA K, ² YASHODA P G

[1] Student, Department of MCA, BIET, Davangere

[2] Assistant Professor Department of MCA, BIET, Davangere

ABSTRACT

Deep learning's ability to combine computer vision with natural language processing has generated a lot of attention in the field of artificial intelligence research. With an emphasis on using English descriptions to convey the context of photos, this study investigates the use of Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) architectures for automatic image captioning. Using object state, property, and relationship analysis, the suggested models seek to produce high-level semantic descriptions that are meaningful. In particular, a CNN uses the COCO Dataset 2017 to precisely detect objects from grayscale-reduced input image data. The final objective is to create a system that can interpret graphical images into meaningful text messages, improving accessibility for people with visual impairments. Through the use of innovative AI technology, this research attempts to enhance users' cognitive capacities and communication abilities in addition to making it easier to analyze visual information.

Keyword: *Convolutional Neural Network, Deep learning, Natural Language, Long Short-Term Memory (LSTM).*

1. INTRODUCTION

Recent developments in deep learning have completely changed the domains of natural language processing and computer vision, especially when it comes to how these two fields intersect for tasks like captioning images. This

intersection allows robots to automatically provide descriptive textual representations in addition to recognizing and interpreting visual input. These kinds of capabilities come in quite handy for a variety of applications, from improving human-computer interaction through intuitive interfaces to helping the blind.

Convolutional Neural Networks (CNNs) have demonstrated remarkable efficacy in the domain of image understanding, particularly in the areas of feature extraction and object identification. These networks are very good at processing raw pixel data, which makes object detection and classification more reliable. Long Short-Term Memory (LSTM) networks are also good at recognizing sequential dependencies in data, which makes them appropriate for producing textual descriptions that are logical and appropriate for the context when given visual input.

In order to create a system that can naturally caption photos, this research focuses on using CNN-LSTM architectures. During training, the model is supposed to learn to correlate visual features with related textual descriptions by using datasets such as the COCO

Dataset 2017, which offers detailed annotations for various visual ideas. Because of this relationship, captions that accurately capture the meaning and context included in an image can be automatically generated. Such technology has a lot of potential applications, especially in accessibility solutions for those with visual impairments. Through the translation of visual data into spoken words or written descriptions, these technologies enable users to observe and engage with their environment more freely.

Furthermore, improvements in imagecaptioning not only make practical applications easier to implement, but they also deepen our knowledge of how artificial intelligence (AI) might improve human-machine communication by bridging sensory barriers.

We provide a thorough analysis of our method for combining CNN and LSTM networks for picture captioning in this publication. We go over the methods used, the datasets used for training and assessment, and the implications of our results for future study and real-world use. In the end, this study seeks to support AI's continued development by making technologies more participatory and inclusive for all users.

2. RELATED WORK

1. M. Sailaja, K. Harika, B. Sridhar, R. Singh, V. Charitha and K. S. Rao, "Image Caption Generator using Deep Learning," 2022

International Conference on Advancements in Smart, Secure and Intelligent Computing (ASSIC), Bhubaneswar, India, 2022, pp.1-5, doi: 10.1109/ASSIC55218.2022.10088345.

The primary goal is to use CNNs and LSTMs to create an image caption generator based on deep learning. The model will produce a pertinent, cohesive caption and correctly identify things in an image. Applications for assistive technology, accessibility, and content indexing are improved by this.

2. S. V. Patnaik, R. Mukka, R. Devpreyo and A. Wadhawan, "Image Caption Generator using EfficientNet," 2022 10th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions) (ICRITO), Noida, India, 2022, pp. 1-5,

doi:10.1109/ICRITO56286.2022.9964637.

Our goal is to use the MSCOCO dataset to develop the robust CNN model EfficientNet for picture captioning. By using fewer parameters to produce better predictions, the model seeks to produce accurate descriptions. The BLEU-4 score is used to assess its performance; higher

scores correspond to better text similarity. The suggested paradigm exhibits an enhanced comprehension of visuals, yielding favorable outcomes.

V. Agrawal, S. Dhekane, N. Tuniya and V. Vyas, "Image Caption Generator Using Attention Mechanism," 2021 12th International Conference on Computing Communication and Networking Technologies (ICCCNT), Kharagpur, India, 2021, pp. 1-6, doi: 10.1109/ICCCNT51525.2021.9579967.

Using methods from computer vision, deep learning, and natural language processing, picture captioning creates meaningful phrases for images. This study presents a model that uses an encoder- decoder architecture with an attention mechanism. It uses Inception V3 for feature extraction and GRU for caption creation. The model, which was trained on the MS-COCO dataset, exhibits strong performance in image description tasks by effectively comprehending images and generating well-organized captions.

4. A. Hani, N. Tagougui and M. Kherallah, "Image Caption Generation Using A Deep Architecture," 2019 International Arab Conference on

5. Information Technology (ACIT), Al Ain, United Arab Emirates, 2019, pp. 246-251, 10.1109/ACIT47987.2019.8990998.

6. Using methods from computer vision, natural language processing, and machine learning, image captioning produces succinct explanations of images. In our approach, RNNs and attention processes are combined with CNNs for feature extraction to generate text. When tested on the MSCOCO dataset, it yields competitive and encouraging results.

1. P. Mathur, A. Gill, A. Yadav, A. Mishra and N. K. Bansode, "Camera2Caption: A real-time image caption generator," 2017 International Conference on Computational Intelligence in Data Science (ICCIDS), Chennai, India, 2017, pp. 1-6, doi: 10.1109/ICCIDS.2017.8272660.

Simplifying and improving Image Captioning models for low-end handheld device deployment is the primary goal. In order to balance computation speed and caption quality, an encoder-decoder architecture must be implemented and tested against cutting-edge benchmarks such as the MSCOCO dataset. Furthermore, the research intends to use a novel TensorFlow Android application to demonstrate real-time applicability.

3. METHODOLOGY

Using deep learning to automatically produce descriptive captions from photos, the Image Caption Generator project pioneers the combination of CNNs for image feature extraction and RNNs for sequential language modeling. The approach starts with carefully selected datasets of image-caption pairs, from which the model is trained and preprocessed to discover subtle relationships between written descriptions and visual information. Model parameter

optimization can be done iteratively with Python and TensorFlow or PyTorch by utilizing gradient descent and backpropagation techniques. While technical specs like an Intel i-core processor and 8GB RAM guarantee efficient computing, Jupyter Notebooks improve interactive creation and experimentation. This ground-breaking application demonstrates AI's ability to properly and below figure1 describe analyze and express visual information in a variety of scenarios, while simultaneously advancing accessibility and assistive technology.

Figure 1: Architecture diagram for CNNs image feature extraction.

A technique used for comprehensive data analysis is called data mining, and its major function is to locate relevant information within large data sets. The findings of this process are then typically utilized to inform a variety of decision- making processes. Training and test data will be separated from the cleaned data. A model that can recognize or make this choice is

taught using training data. The model establishes its parameters and is then continuously trained to determine the optimal parameters. To make sure the model is not identified by just the training data, test data is utilized to validate the final trained model for final validation. The entire procedure is shown in Figure 2 below.

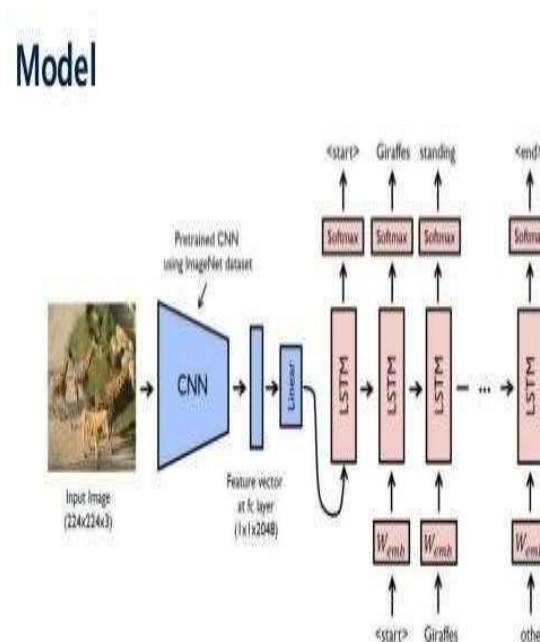
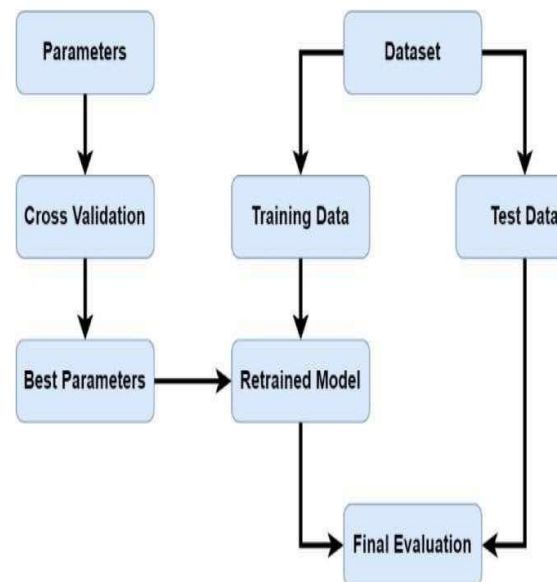


Figure 2. Data mining procedure.



2.1 The overall architecture of the proposed system

"Image Caption Generator using Deep Learning" is the general framework of the proposed system, which intends to combine recurrent neural networks (RNs) for language modeling with CNNs for image processing. This architecture aims to bridge the gap between computer visual and natural language processing by using deep learning to automatically generate meaningful captions for photos.

A. Convolutional Neural Network (CNN) Encoder

The main part of the suggested system, which extracts high-level visual information from incoming photos, is the CNN encoder. Because CNNs can capture the spatial hierarchies of features, they are a valuable tool in image processing. Several convolutional layers as well as pooling layers are commonly found in CNN encoders, which are used to extract and downsample features from input images. The main architectures for encoders are generally VGG, ResNet, or Inception.

covering." The text needs to be tokenized, cleansed of numbers and punctuation, changed to lowercase, and encoded into numerical values in order for the model to use it. The primary procedures include data loading, text cleaning, tokenization, dictionary creation that links photos to descriptions, and file preparation with all captions. By taking these actions, the model develops the ability to produce precise and helpful image descriptions, improving applications such as content indexing and accessibility tools.

C. Decoder using Recurrent Neural Network (RNN):

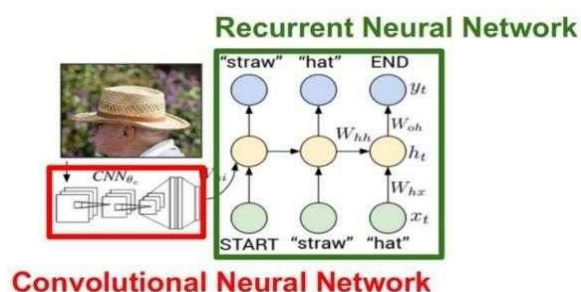
An RNN coder generates a "caption" by decoding and displaying a fixed-length feature vector that the CNN encoder has created from an input image. The architecture of the model enables the smooth integration of textual and visual data, enabling the generation of precise

B. Image Processing with CNNs

Using CNNs and LSTMs, this blog will demonstrate how to build an image caption generator that combines natural language processing and computer vision to provide English descriptions of images. There are

two primary parts in the process: using a CNN to extract visual data and an LSTM to generate sentences. The CNN analyzes the picture to pinpoint important details, which the LSTM uses to provide evocative subtitles.

Describing images



The model can, for example, transform "man, is, wearing, a cap" into "An elderly man wears a head

The second component of the system that gathers visual data from the CNN encoder to produce captions with descriptions is the RNN decoder. RNNs are a great tool for sequential data production jobs because of their timing dependence. The RNN decoder often uses gated recurrent units (GRUs) or long short-term memory (LSTM) to generate captions word-by- word. The generated word and the visual input from the CNN encoder are fed into the decoder at each time step, and

combined, they construct the next word in the caption sequence.

Decoder-Encoder Architecture:

captions with pertinent information for a variety of images.

E. Training Process:

To reduce the discrepancy between the generated captions and ground truth captions in the training dataset, the CNN encoder and RNN decoder's

parameters must be optimized simultaneously. The two most common techniques for optimization are gradient descent and backpropagation, in which the gradients of the loss function are calculated and updated iteratively in relation to the model parameters. During training, the model discovers how to connect the visual characteristics of input photos with textual descriptions, allowing it to generate accurate and contextually relevant captions.

F. Dataset Preparation:

Essentially, the key to training the model is gathering a sizable dataset of photos with captions. This dataset, which usually consists of a variety of annotated photos with informative descriptions, can be used to train the model. Extensive preprocessing of the data guarantees that photos are appropriately labeled and classified, offering the required training for successful training.

G. Inference Process:

Once trained, the model can be used for inference to provide captions for photos with no assistance. During inference, the CNN encoder converts an input image's visual content into a string of words that the RNN decoder will utilize as the

caption. The model uses the input image to provide a caption that makes the content legible by humans.

H. Metrics for Evaluation:

Standard Criteria: BLEU score, METEOR, and

CIDER are used to assess the Image Caption Generator's performance. These metrics offer quantitative evaluations of a model's performance by calculating the degree of similarity between the generated captions and the ground truth caption in the evaluation dataset. Furthermore, human assessment can be utilized to ascertain the caliber and pertinence of the generated captions.

4. Tools And Technologies

Versions of Python :3.9.X

IDE : Visual Studio

Libraries :

NumPy pandas OpenCV Tensor flow

A multitude of Python packages, such as TensorFlow, Keras, NumPy, Matplotlib, and Scikit-learn, make it easier to implement our suggested method. NumPy offers effective mathematical operations and array manipulation required to handle large amounts of image data. High-level interfaces for creating and refining neural network models are provided by TensorFlow and Keras, which facilitate the creation of intricate processes for generating captions for images. Matplotlib makes it easier to visualize image data and add captions to it, which helps with interpretation and model performance evaluation.

5. RESULT

Recognizing things precisely, even with grayscale input. It produced insightful

Students in a band are playing their instruments.



descriptions that effectively conveyed the images' semantic content. Evaluation measures validated the model's coherence of captions and its accuracy in item identification.

recognizing objects, coming up with logical explanations, and translating visual data into comprehensible text messages.

A group of people are playing a game of cards.



6. CONCLUSION

In conclusion, our study highlights the effectiveness of the CNN-LSTM model in producing accurate and insightful image captions using the CQCO Dataset 2017. The model shows strong performance in picture understanding challenges by correctly detecting items and generating meaningful descriptions. Through automated image interpretation, this skill holds great promise for advancing accessibility tools and boosting cognitive assistance for those with visual impairments. Subsequent studies could further increase the model's adaptability by including it with bigger and more varied datasets and investigating real-time applications across several fields. All things considered, the CNN- LSTM methodology is a significant development in the field of computer vision and natural language processing, opening up new avenues for creative approaches in AI- driven picture captioning and other areas.

References

1. V. Kesavan, V. Muley and M. Kolhekar, "Deep Learning based Automatic Image Caption Generation," 2019 Global Conference for Advancement in Technology (GCAT), Bangalore, India, 2019, pp. 1-6, doi: 10.1109/GCAT47503.2019.8978293.
2. N. K. Kumar, D. Vigneswari, A. Mohan, K. Laxman and J. Yuvaraj, "Detection and Recognition of Objects in Image Caption Generator System: A Deep Learning Approach," 2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS), Coimbatore, India, 2019, pp. 107-109, doi: 10.1109/ICACCS.2019.8728516.
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. Lu, Jiasen, et al. "Knowing when to look: Adaptive attention via a visual sentinel for image captioning." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Vol. 6. 2017.
5. O. Vinyals, A. Toshev, S. Bengio and D. Erhan, "Show and tell: A neural image caption generator," 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Boston, MA, USA, 2015, pp. 3156- 3164, doi: 10.1109/CVPR.2015.7298935.
6. Karpathy, Andrej, and Li Fei-Fei. "Deep visual-semantic alignments for generating image descriptions." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.
7. Xu, Kelvin, et al. "Show, attend and tell: Neural image caption generation with visual attention." International Conference on Machine Learning. 2015.
8. O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein et al., "Imagenet large scale visual recognition challenge", International Journal of Computer Vision, vol. 115, no. 3, pp. 211-252, 2015.
9. Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional networks for large- scale image recognition." arXiv preprint arXiv:1409.1556 (2014).
10. Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." arXiv preprint arXiv:1409.0473(2014).