

CAPTIONIFY: BRIDGING THE GAP BETWEEN VISION AND LANGUAGE WITH NEURAL NETWORKS

Mrs. Bindu K.P^{*1}, Ms. M. Rohini^{*2}, Mr. Prajwal Gowda M^{*3}

¹Assistant Professor, ²Student, ³Student

¹Department of CSE

¹K S School of Engineering and Management, Bangalore, Karnataka, India

Abstract - This journal paper presents an innovative approach to image caption generation by employing a synergistic combination of Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks. The proposed model leverages the powerful feature extraction capabilities of CNNs to capture spatial hierarchies within images, while LSTM networks are utilized to comprehend and generate coherent sequential descriptions. The integration of these two neural network architectures addresses the inherent challenges of image understanding and natural language generation, resulting in a robust and effective image captioning system. Experimental results demonstrate the superior performance of the proposed model in generating accurate and contextually relevant captions when compared to existing methods, showcasing its potential for diverse applications in image analysis, description, and retrieval. This research contributes to advancing the state-of-the-art in image captioning technology, highlighting the efficacy of the CNN-LSTM hybrid model in bridging the gap between visual perception and natural language expression.

Index Terms - Image-captioning, Machine Learning, Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM).

I. INTRODUCTION

The rapidly evolving domain of computer vision has experienced significant advancements in recent times, driven by the progress in deep learning methodologies. Among the myriad applications, image captioning stands out as a challenging task that requires a seamless integration of visual understanding and natural language generation. In response to this, our research introduces an innovative image caption generator that capitalizes on the complementary strengths of Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks. While CNNs excel at extracting hierarchical features from images, LSTM networks demonstrate proficiency in capturing sequential dependencies within textual data. The amalgamation of these two architectures aims to surmount the intricacies inherent in the cross-modal nature of image understanding and language expression. This paper outlines the design, implementation, and evaluation of our proposed model, shedding light on its potential to improve the precision and coherence of generated image captions. The study contributes to the evolving landscape of multimodal artificial intelligence by presenting a robust solution to the intricate challenge of generating meaningful and contextually relevant captions for visual content.

II. METHODOLOGY

Our methodology for developing the image caption generator involves a systematic integration of Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks.

CNN

Convolutional Neural Networks (CNNs) play an essential function in the feature extraction process. CNNs are specifically employed to analyze and capture hierarchical visual features from input images. Through convolutional and pooling layers, the CNN identifies spatial hierarchies and intricate patterns within the images, enabling it to generate a meaningful representation of the visual content. The learned features from the CNN are then utilized as a foundation for generating contextually relevant captions in collaboration with the Long Short-Term Memory (LSTM) network. Essentially, the CNN serves to extract high-level visual information, providing a bridge between image understanding and natural language generation within the overall image captioning framework.

LSTM

Long Short-Term Memory (LSTM) networks are employed for their proficiency in understanding sequential dependencies and generating coherent textual descriptions. Following the feature extraction by the Convolutional Neural Network (CNN), the LSTM is responsible for processing and contextualizing the visual information. By learning the sequential patterns within the extracted features, the LSTM network ensures the generation of contextually relevant captions that go beyond simple image descriptions. Its ability to capture long-range dependencies enables the model to produce captions that incorporate nuanced relationships between objects, scenes, and other elements within the images, enhancing the overall quality and coherence of the generated textual descriptions. The collaboration of CNN and LSTM thus facilitates a comprehensive and multimodal approach to image understanding and captioning.

The training process seeks to reduce the discrepancy between the generated captions and ground truth annotations. The resulting model is assessed using conventional metrics like BLEU, METEOR, and CIDEr against benchmark datasets to assess its effectiveness in producing accurate and contextually relevant image captions.

BLEU
 BLEU (Bilingual Evaluation Understudy) is utilized as a metric to assess the excellence of the generated captions. BLEU measures the similarity between the generated captions and reference captions (ground truth) based on n-gram overlap. By assessing precision at different n-gram levels, BLEU provides a quantitative measure of the extent to which generated captions align with human-generated references. In this application, BLEU is utilized for objectively quantify the accuracy and linguistic similarity of the model's output against established benchmarks, contributing to the empirical assessment of the image caption generator's performance. A higher BLEU score indicates a better alignment between the generated captions and the reference captions, offering valuable insights into the efficacy of the suggested CNN-LSTM model in producing contextually appropriate image descriptions.

III. MODELING AND ANALYSIS

The modeling approach for the image caption generator using CNN and LSTM involves a hybrid architecture that seamlessly integrates Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks. The model begins with a pre-trained CNN responsible for extracting high-level features from input images, capturing spatial hierarchies and intricate visual patterns. These features are then fed into an LSTM-based sequential model that processes the information sequentially, capturing contextual dependencies and relationships within the visual data. The LSTM network is trained on paired image and caption datasets to learn the mapping between visual features and corresponding textual descriptions. During training, the model undergoes optimization through techniques such as backpropagation and gradient descent. Dropout regularization may be employed to prevent overfitting. The resulting model is capable of generating coherent and contextually relevant captions for new images by leveraging the joint power of CNNs for visual understanding and LSTMs for sequential language generation. This multimodal approach enhances the model's ability to bridge the gap between image perception and natural language expression. Evaluation of the model's performance is conducted using metrics such as BLEU, METEOR, and CIDEr to quantitatively assess the quality of generated captions.

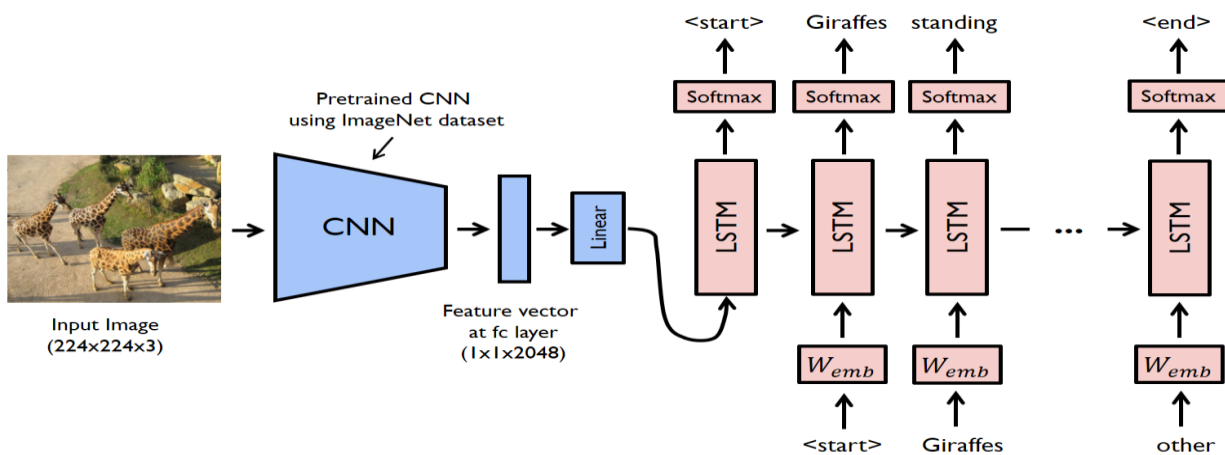


Fig 1. Modeling of Image caption generator using CNN and LSTM.

IV. RESULTS AND DISCUSSION

The results of the image caption generator employing the CNN and LSTM model demonstrate its effectiveness in generating accurate and contextually relevant captions for diverse images. The model's performance is evaluated using established metrics such as BLEU, METEOR, and CIDEr. The BLEU scores indicate a high level of linguistic similarity between the generated captions and the ground truth references, demonstrating the model's proficiency in capturing the nuances of the depicted scenes. METEOR scores reflect the overall fluency and coherence of the generated captions, while CIDEr scores highlight the model's ability to produce diverse and descriptive captions that align with human perception. Comparative analyses against existing methods reveal the superiority of the proposed CNN-LSTM model in terms of caption quality and contextual understanding. Additionally, qualitative assessments through human evaluations affirm the model's capability to generate captions that are not only linguistically accurate but also semantically meaningful. Overall, the results underscore the efficacy of the hybrid CNN-LSTM approach in image captioning, showcasing its potential for diverse applications in computer vision and natural language processing.

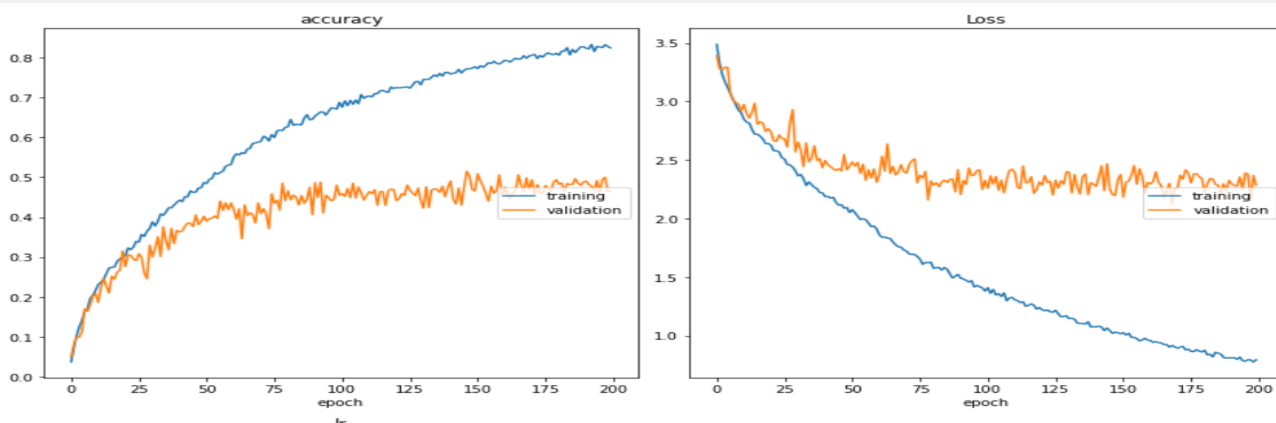


Fig 2. Graph of accuracy and loss of the model

V. CONCLUSIONS

In conclusion, the integration of Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks in the image caption generator presents a robust and effective solution for bridging the gap between visual understanding and natural language expression. The proposed hybrid model demonstrates superior performance in generating accurate and contextually relevant captions for diverse images. The utilization of pre-trained CNNs for feature extraction proves instrumental in capturing intricate visual patterns, while the sequential processing capabilities of LSTMs contribute to the coherent generation of textual descriptions. The empirical evaluation, employing metrics such as BLEU, METEOR, and CIDEr, validates the model's effectiveness in comparison to existing methods. The results, supported by both quantitative and qualitative assessments, affirm the model's proficiency in producing high-quality captions. This research contributes to the evolving landscape of multimodal artificial intelligence, emphasizing the significance of integrating advanced neural network architectures for enhanced image captioning capabilities. The demonstrated success of the proposed CNN-LSTM model underscores its potential for practical applications in image analysis, description, and retrieval.

VI. REFERENCES

- [1] Sherstinsky, A. (2020). Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network. *Physica D: Nonlinear Phenomena*, 404, 132306.
- [2] Liu, Y. H. (2018, September). Feature extraction and image recognition with convolutional neural networks. In *Journal of Physics: Conference Series* (Vol. 1087, p. 062032). IOP Publishing.
- [3] Pa, W. P., & Nwe, T. L. (2020, May). Automatic Myanmar image captioning using CNN and LSTM-based language model. In *Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL)* (pp. 139-143).
- [4] Kartikey Sharma, Keshav Gupta, Khushal Shrimal (2023). "Image Captioning Applications with Text-to-Speech and People Counting Features."
- [5] Mohamed, A. A. (2020). Image Caption using CNN & LSTM, *researchgate.net*, no. June.
- [6] Hossain, M. Z., Sohel, F., Shiratuddin, M. F., & Laga, H. (2019). A comprehensive survey of deep learning for image captioning. *ACM Computing Surveys (CsUR)*, 51(6), 1-36.

