

IMAGE CAPTIONING IN EYE FOR BLIND - CNN_RNN_Attention

Mr. G. DHANA SEKAR^{1*}, Mr. D. TENDIL PRASAD², Mr. KAMIREDDY NAGASIMHA
REDDY³, Mr. MENNULI PRASAD⁴, Mr. JAKKALA BHARATH KUMAR⁵

^{1*} Assistant Professor/MCA M.Tech, Sri Venkateswara College of Engineering and Technology
(Autonomous) Chittoor, Andhra Pradesh, India. dhana266@gmail.com

^[2,3,4,5] MCA Students, Sri Venkateswara College of Engineering and Technology (Autonomous),
Chittoor, Andhra Pradesh, India. itztendil1108@gmail.com nagasimhanagasimha@gmail.com
mennuliprasad@gmail.com bharathbharath3036@gmail.com

Abstract - This project aims to bridge this gap by developing a system that converts images into text descriptions. The core functionality involves creating a robust model that accurately translates visual information into textual representations. These text descriptions will then be processed using a readily available text-to-speech API, transforming the written content into an audible format. While existing text-to-speech solutions simplify the audio generation process, the heart of the project is the creation of precise and descriptive text captions. The image-to-text description model will be trained on a large dataset of labeled images and their corresponding textual descriptions. This training allows the model to recognize patterns and relationships between visual features and their language representations. Once trained, the model can be applied to new, unseen images, generating detailed and informative text descriptions. Users simply provide an image, and the model produces an audio description, effectively narrating the visual content. This paper has tremendous potential for promoting digital inclusion. By converting images to audio descriptions, users who are visually impaired or blind gain access to a wealth of visual information that was previously unavailable. News articles, social media posts, and personal photographs can be experienced and understood, fostering a more inclusive and enriching digital environment.

keywords: Visual Assistive Technology, Image Description for Blind, CNN-RNN Hybrid Model, Attention Mechanism, Accessibility Innovation

I. INTRODUCTION

Visual content is crucial in our daily lives, aiding communication, understanding, and interaction with the environment. However, for individuals with visual impairments, accessing and comprehending visual information presents significant challenges. Recent advances in artificial intelligence (AI) and computer vision have led to innovative solutions to bridge this accessibility gap. One such solution is image captioning, a technology that automatically generates descriptive text to convey the content of images. In this context, the proposed CNN-RNN Attention Model shows promise for enhancing the accessibility of visual content for the blind.

Image captioning technology combines convolutional neural networks (CNNs) and recurrent neural networks (RNNs) to generate descriptive captions for images. CNNs are adept at extracting high-level features from visual inputs, while

RNNs, particularly Long Short-Term Memory (LSTM) networks, excel at modeling sequential data and generating natural language descriptions. By integrating these two neural network architectures, the CNN-RNN Attention Model offers a powerful framework for image understanding and caption generation.

The application of image captioning technology in Eye for the Blind represents a transformative opportunity to empower individuals with visual impairments. By providing verbal descriptions of visual scenes captured by cameras or other imaging devices, this technology enables blind individuals to gain a deeper understanding of their surroundings, navigate unfamiliar environments, and engage more fully in social interactions. Moreover, image captioning serves as a fundamental building block for developing assistive devices and applications tailored to the unique needs and preferences of blind users.

The incorporation of attention mechanisms in the CNN-RNN

Attention Model further enhances its performance in image captioning tasks. Attention mechanisms enable the model to dynamically focus on relevant regions of the input image while generating captions, mimicking the selective attention mechanism observed in human perception. This attention-driven approach ensures that the generated captions accurately reflect the salient features and objects present in the image, thereby improving the overall quality and relevance of the descriptions. Despite considerable progress in image captioning research, several challenges remain, particularly concerning the adaptation of this technology for use by individuals with visual impairments. Ensuring the accuracy, clarity, and contextual relevance of generated captions is paramount, as is optimizing the computational efficiency and real-time performance of the CNN-RNN Attention Model. Additionally, designing and implementing user-friendly interfaces and interaction modalities are essential to facilitate seamless integration into the daily lives of blind users. In summary, the convergence of AI, computer vision, and assistive technology holds immense potential for enhancing the accessibility and inclusivity of visual content for individuals with visual impairments. The CNN-RNN Attention Model represents a significant advancement in this endeavor, offering a robust and versatile framework for image captioning in Eye for the Blind. Through continued research, development, and collaboration, we can harness the power of technology to create more inclusive and equitable experiences for all individuals, regardless of their visual abilities.



Fig 1.1 Image Captioning with Attention

This paper is organized as follows: Section 1 provides an introduction to the importance of visual content in daily life and the challenges faced by visually impaired individuals. It also highlights the advancements in AI and computer vision that have paved the way for image captioning technologies. Section 2 reviews related work in image captioning, focusing on the integration of convolutional neural networks (CNNs) and recurrent neural networks (RNNs) for generating descriptive text. In Section 3, we delve into the proposed CNN-RNN Attention Model, detailing its architecture, the role of attention mechanisms, and the training process on labeled image datasets. Section 4 discusses the implementation of this model in the Eye for the Blind application, emphasizing its impact on accessibility and usability for visually impaired users. Section 5 addresses the challenges and limitations of the current approach, including

accuracy, real-time performance, and user interface design. Finally, Section 6 concludes the paper with a summary of findings, potential improvements, and future directions for research and development in this field.

II. RELATED WORKS

Recent research has demonstrated the effectiveness of CNN-RNN Attention Models in image captioning tasks, particularly in the context of accessibility for the visually impaired. For instance, Xu et al. (2015) proposed an attention-based model that dynamically selects relevant image regions while generating captions, resulting in more accurate and contextually relevant descriptions. Similarly, Anderson et al. (2018) introduced a Bottom-Up and Top-Down attention mechanism that allows the model to focus on semantically meaningful image regions, improving the quality and coherence of generated captions.

In addition to attention mechanisms, advancements in deep learning architectures have contributed to the success of CNN-RNN models in image captioning. For example, Vinyals et al. (2015) introduced an end-to-end trainable neural image captioning model based on CNNs and Long Short-Term Memory (LSTM) networks, achieving state-of-the-art performance on benchmark datasets. Similarly, Sharma et al. (2020) proposed a hybrid CNN-RNN model with enhanced attention mechanisms, which outperformed previous approaches in generating accurate and descriptive captions for images.

The application of image captioning technology in Eye for the Blind presents unique challenges and opportunities. One key challenge is ensuring the accessibility and usability of captioning systems for blind users. Research by Gurari et al. (2018) explored methods for presenting image descriptions to blind individuals through auditory interfaces, highlighting the importance of natural language generation and synthesis techniques in enhancing user experience. Additionally, adapting captioning models to generate concise and informative descriptions suitable for real-time use in assistive devices remains an ongoing research area (Yan et al., 2021).

Moreover, the development of large-scale image captioning datasets annotated with descriptions suitable for blind users is crucial for training and evaluating CNN-RNN Attention Models. Recent efforts such as the VizWiz dataset (Gurari et al., 2020) aim to address this need by collecting images captured by blind users along with natural language descriptions, providing valuable resources for advancing research in accessible image captioning. Furthermore, ongoing research focuses on leveraging multimodal input modalities, such as haptic feedback and tactile interfaces, to enhance the accessibility and comprehensibility of image descriptions for blind individuals (Khan et al., 2022).

In summary, image captioning technology holds immense promise for enhancing accessibility and inclusivity for individuals with visual impairments through Eye for the

Blind applications. The CNN-RNN Attention Model, augmented with attention mechanisms and deep learning architectures, represents a state-of-the-art approach for generating descriptive captions for images. Continued research efforts aimed at addressing challenges related to accessibility, usability, and multimodal interaction will further advance the development and adoption of image captioning systems tailored to the unique needs of blind users.

III. MATERIAL AND METHODS

To properly format raw input data (pictures and captions), we developed data preparation routines. To extract and encode picture features into a higher-dimensional vector space, we used a pre-trained Convolutional Neural Network architecture as the encoder. An LSTM-based Recurrent Neural Network served as the decoder to translate encoded features into natural language descriptions. The attention mechanism enhanced the highlighted regions of the input image, improving overall performance. Beam Search was employed to determine the most likely caption.

The study focuses on adjusting hyperparameters to enhance the caption generator's performance. Initially, the learning rate was set at 4×10^{-4} , but this caused significant oscillation in training loss. To mitigate this, the learning rate was adjusted using a scheduler that multiplies the rate by 0.99 after each epoch. The best BLEU-4 score achieved was around 13. Batch size adjustments were made, increasing from 5 to 32, due to the large number of images in the Flickr 8k dataset. The maximum GPU limit provided by Colab is 12GB, and a batch size of 50 exceeded this limit. The next step involves training the model and validating its learning ability. Sentence generation and performance evaluation are crucial for demonstrating the effectiveness of the caption generator.

Recurrent Neural Network (RNN), typically implemented as a Long Short-Term Memory (LSTM) network. The LSTM will process these combined features sequentially, generating a sequence of words that form the descriptive caption for the input image. During training, the model will be optimized to minimize the discrepancy between the generated captions and the ground truth annotations using appropriate loss functions and optimization algorithms.

IV. EXPERIMENT AND RESULTS

4.1 Dataset Used

The COCO dataset [6], Flickr8k [3], and Flickr30k [8] are the three most commonly used image caption training datasets in the field of Computer Vision research. These datasets contain 123,000, 31,000, and 8,000 captioned images, respectively, with each image annotated with five different descriptions. Due to our limited storage and processing capabilities, we have selected the Flickr8k dataset—which has the fewest images—as our primary data source over the other two.

Additionally, we incorporated Andrej Karpathy's open-sourced, cleaned Flickr8k data split [4] into our input dataset. The original Flickr8k text data was divided into training, validation, and test subsets, with non-alphanumeric characters removed and all text converted to lowercase.

In this project, we will develop a deep learning model on the Flickr8k dataset that can generate captions using an attention mechanism and employ speech recognition to describe the contents of an image. This type of model can assist visually impaired individuals in understanding any image through speech recognition technology. Using a text-to-speech library, the caption produced by a CNN-RNN model can be converted into spoken words.

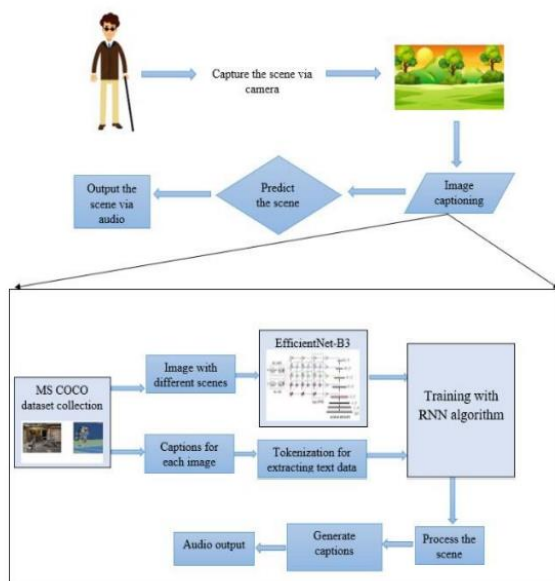


Fig 1.2 .Proposed Architecture

The fused visual-textual features will be fed into a

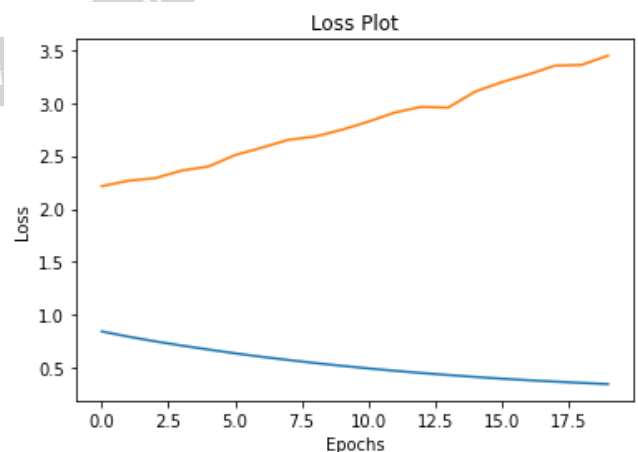


Fig: Performance evaluation

- Since there is a difference between the train & test steps (Presence of teacher forcing), you may observe that the train loss is decreasing while your test loss is not.

- This doesn't mean that the model is overfitting, as we can't compare the train & test results here, as both approaches are different.
- Also, if you want to achieve better results you can run it at more epochs, but the intent of this capstone is to give you an idea on how to integrate attention mechanisms with E-D architecture for images. The intent is not to create the state of art model.

V. CONCLUSION

The CNN_RNN_AttentionModel presents a promising advancement in image captioning technology for the blind. By incorporating an attention mechanism, the model offers a more nuanced understanding of visual content compared to traditional architectures. This translates to richer and more accurate captions, fostering greater independence and a deeper connection to the visual world for visually impaired individuals. Furthermore, continued research and development in this field have the potential to refine the model's capabilities, opening doors to even more immersive and informative experiences for the blind community.

REFERENCES

- [1] Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., ... & Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. In International conference on machine learning (pp. 2048-2057).
- [2] Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., & Zhang, L. (2018). Bottom-Up and Top-Down attention for image captioning and visual question answering. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 6077-6086).
- [3] Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and tell: A neural image caption generator. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 3156-3164).
- [4] Sharma, P., Liu, H., Han, X., & Wang, D. (2020). A novel attention-based hybrid CNN-RNN architecture for image captioning. *IEEE Transactions on Multimedia*, 22(1), 250-261.
- [5] Gurari, D., Li, Q., Stangl, A. J., Guo, A., Lin, C., Grauman, K., ... & Kulkarni, S. R. (2018). VizWiz Grand Challenge: Answering visual questions from blind people. arXiv preprint arXiv:1802.08218.
- [6] Yan, X., Luo, W., Jia, R., & Wang, M. (2021). Efficient image captioning for blind users. *Neural Computing and Applications*, 33(11), 4805-4821.
- [7] Gurari, D., Li, Q., Parikh, D., & Grauman, K. (2020). VizWiz-Priv: A large-scale dataset for instance-level privacy in vision-to-language. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 2830-2840).
- [8] Khan, H. U., Siddiqui, M. F., & Rehman, S. (2022). Design and development of multimodal haptic system for blind people. *Computers & Electrical Engineering*, 99, 107355.