

# AI-based Image Captioning and Scene Description

Puneet Kaur<sup>1</sup>, Astha Srivastava<sup>2</sup>, Priya Kumari<sup>3</sup>, Sonu Kumar Saw<sup>4</sup>, Lakhani Singh<sup>5</sup>  
<sup>1,2,3,4,5</sup> Department of Computer Science Chandigarh University Mohali, India

**Abstract:** *In recent improvements of image captioning and scene description based on AI, there has been great enhancement in the areas of accessibility for visually impaired users and management of contents. Recent works focus on new approaches that combine deep learning models, transformers, and multimodal techniques for high-quality, context-sensitive image description. Exploiting these emerging technologies, our research introduces a newly designed device that targets visually impaired users. We are proposing in this paper a machine learning, large language modeling, and natural language processing device that is able to voice out the details captured by a camera regarding an object and provide a description. Our device incorporates many state-of-the-art techniques proposed in recent studies, such as semantic and visual attention mechanisms, and delivers even better accuracy with more contextual relevance in its descriptions. This approach has contributed not only to developing the field of assistive technologies but also to the greater goal of making visual information more accessible and understandable by the blind.*

**Keywords:** *AI-based image captioning, scene description, accessibility, visually impaired, machine learning, large language models (LLM), natural language processing (NLP), assistive technologies, deep learning, multimodal techniques.*

## I. INTRODUCTION

AI-based image captioning and scene description have achieved remarkable development powered by the emergence of ML and DL technologies. Transformers, attention mechanisms, and LLMs are among the key approaches that have emphatically enhanced the accuracy of automatic captions of images, in regards to accessibility for the visually impaired.

Zhang et al. (2023) identified that Transformers have been helpful in the generation of scene graphs, a method for representing both spatial and semantic among objects of an image. This has proved to be the cornerstone in generating more meaningful and accurate descriptions of images that make more sense in helping visually impaired users understand their environment. According to their view by Wang et al. (2023), deep learning models had great ability in

generating detailed captions due to the handling of complex visual data, enhancing the ability of assistive technologies by a great margin. Attention mechanisms—visual and semantic attention—play the most imperative roles in creating coherent and contextually relevant caption generation. Huang et al. showcased that the attention mechanism within models in image captioning enhances precision by shifting focus on the most important aspects of an image. These models help refine generated captions so that they become more informative and relevant for visually impaired individuals [3]. However, despite this, they tend to face problems in handling complex scenes with many objects having complex relationships. Wang et al. (2023) attempted to solve this problem using graph neural networks that enhance scene understanding; these approaches are usually computationally expensive and hard to apply in real-world applications [4].

Another major problem with AI-based image captioning is in the real-time processing of images. While textual and visual multimodal approaches enhance the performance of image descriptions, most of these models are computationally resource-intensive, hence challenging the possibility of real-time feedback. Kumar et al. (2023) explored these multimodal methods but noted that real-time caption generation was one of the limitations, especially for assistive devices meant for visually impaired users [5]. Xu et al. (2023) find that real-time feedback is needed on such a device, and instantaneous response becomes critical in order to communicate with the user interactively with success [6].

Addressing these challenges, our research proposes a wearable device by visually impaired individuals for real-time object scanning, with detailed descriptions through voice output. The incorporation of machine learning, large language models, and natural language processing enhances object recognition and description accuracies within this device. Unlike previous approaches, our device embeds real-time processing using advanced attention mechanisms that allow immediate and contextually rich information to

the visually impaired user [7][8].

## II. LITERATURE REVIEW

Unlike traditional approaches to multimodal techniques suffering from real-time applicability, our device uses efficient algorithms that balance accuracy with computational demand. This makes the device effective for real-time applications without compromising performance, hence a workable solution to the limitations brought about by existing assistive technologies [9][10].

This is because of the pace in which image captioning and scene description are gaining importance with advancements in ML and DL. Some robust technologies, namely, transformers, attention mechanisms, and LLMs, have contributed a lot in this area by making the accuracy of automatically generated image captions much better for accessibility by visually impaired persons.

Year	Author	Title	Contribution	Limitation	Relevance to Current Study
2023	Zhang, H., Wu, W., & Chen, L.	Scene Graph Generation with Transformers	Introduces transformers to understand object relationships in images, enhancing scene understanding.	Limited to complex scenes; requires a lot of computational resources.	Foundational techniques in understanding object interactions; relevant to improvements on image description systems.
2023	Wang, Q., Li, Z., & Zhang, Y.	Deep Image Captioning: A Review	Reviews a large number of deep learning methods for image captioning and outlines main techniques and trends.	The paper summarizes existing methods but does not really provide new solutions.	Provides overview of deep image captioning methods that will inform improvements in image description systems.
2023	Huang, J., Yang, X., & Zhang, T.	Image Captioning with Visual and Semantic Attention	This is an advanced version of image captioning with a combination of visual and semantic attention.	It may have scalability and performance issues during real-time applications.	For the development of the system which requires detailed descriptions of images and with any context.
2023	Kumar, S., Roy, A., & Singh, M.	Multimodal Approaches to Visual Question Answering and Image Captioning.	Discusses different multimodal approaches using data from images and text towards better question-answering and captioning.	May not discuss the challenges to integrate different modalities.	Give insights into the multimodal approaches that can be applied to improve the accessibility features present in image captioning.

2023	Wang, L., Liu, Y., & Patel, R.	Object Detection and Captioning in Complex	Uses graph neural networks to improve object detection and	High computational requirements and complexity.	Offers advanced techniques for object detection and captioning, applicable
		Scenes with Graph Neural Networks.	captioning in complex scenes.		to complex scene analysis in assistive technologies.
2023	Xu, H., Li, T., & Lee, D.	Generative Models for Image Captioning and Visual Storytelling	Applies generative models to create engaging captions and narratives from images.	Limited to creative applications; may not handle all practical scenarios	Supports the creation of detailed and narrative-driven descriptions, useful for accessibility enhancements.
2023	Brown, J., Patel, N., & Zhang, S.	Contextual Image Captioning Using Large Language Models	This model leverages large language models to produce context-sensitive image descriptions.	Although reliant on the pre-trained model considerably, it may not adapt in each particular context.	Applied directly to writing contextually appropriate descriptions for users with visual impairments.
2023	Lee, R., Kim, M., & Zhao, Y.	End-to-End Learning for Real-Time Image Captioning and Visual Dialogue	Covers end-to-end learning for real-time image captioning and dialogue.	Limited real-time performance, less adaptability.	This area can offer solutions from real-time into interactive applications in accessibility technologies.
2024	Sharma, A., Gupta, P., & Patel, K.	Improving Accessibility for the Visually Impaired with AI-Based Image Description Systems	Focuses on AI systems for generating real-time image descriptions to aid visually impaired individuals.	New approach; still requires practical validation and user testing.	Directly aligns with the goal of improving accessibility through enhanced image description systems.
2024	Johnson, M., Clarke, E., & Wong, D.	Assistive Technologies for the Visually Impaired: Advances and Applications	Covers recent developments in assistive technology for the visually impaired.	Probability-Those latest developments might not have been included that are beyond the review.	Gives a good overview of the assistive technologies which are relevant to the broader task of development of image description systems.

### III. METHODOLOGY

The methodology for developing an AI-based device to assist visually impaired individuals integrates advanced Machine Learning (ML), Deep Learning (DL), and Natural Language Processing (NLP) techniques. This device allows users to scan objects and receive detailed descriptions through earphones, enhancing their interaction with the environment. The methodology leverages pre-trained models and transformer-based architectures, alongside custom techniques to ensure accuracy, relevance, and accessibility for visually impaired users.

#### 1. Data Collection and Pre-processing:

The initial phase involves acquiring large-scale datasets such as MS COCO and Visual Genome, which contain image-caption pairs across a variety of scenarios, ensuring the model can generalize well in real-world environments. Preprocessing will focus on extracting features using Convolutional Neural Networks (CNNs), specifically ResNet and Efficient Net, which decompose each image into key objects and their attributes. Scene segmentation and object detection using models like Faster R-CNN help identify individual objects, while scene graph generation enhances the contextual understanding of object interactions, a process elaborated in [1]. The comprehensive nature of these datasets will enable training for both static and dynamic environments, critical for real-time usability.

#### 2. Model Architecture:

The core architecture integrates transformer networks to handle both visual and textual data. As outlined in [2], transformers efficiently capture long-range dependencies within the scene, ensuring accurate caption generation. A unique idea introduced here is the use of multimodal fusion by combining visual data with temporal data, such as video frames or sequential images, which adds depth to the contextual understanding of object relationships. By leveraging Vision Transformers (VT) for vision tasks and GPT-like models for textual generation, the system can generate high-quality captions.

Additionally, the architecture employs a Graph Neural Network (GNN) to handle complex scenes by modelling object relationships, especially in multi-object environments such as crowded spaces. The GNN refines the system's ability to describe

interactions, offering a more comprehensive view of how objects in a scene relate to one another, as seen in [5].

#### 3. Multimodal Learning and Attention Mechanisms:

The system utilizes multimodal learning, processing both the image and natural language data simultaneously. Attention mechanisms play a critical role by ensuring that the model focuses on the most relevant elements in the image. Visual attention mechanisms allow the model to prioritize essential objects or parts of the scene, while semantic attention ensures the generation of contextually meaningful descriptions, as discussed in [3]. This integration is essential for generating coherent and accurate scene descriptions, especially for visually impaired users who rely on detailed information to navigate their surroundings.

#### 4. Training Process:

The device employs an end-to-end learning framework for training, integrating visual data from CNNs and contextual data from NLP models. This allows the system to generate real-time image descriptions and visual dialogues. The device will also use real-time feedback to adapt to individual user needs, as inspired by [6]. The training will focus on the BLEU, METEOR, and CIDEr metrics to evaluate the accuracy and fluency of the captions generated. A special emphasis will be placed on generating descriptions that are relevant, concise, and adapted to the specific requirements of visually impaired users, as highlighted in [7].

#### 5. Assistive Technology Integration:

Large Language Models (LLMs) and NLP techniques will be used to provide personalized and context-sensitive descriptions, ensuring that users receive accurate, relevant information about their surroundings. These models will be trained to understand both the visual scene and the user's specific needs, allowing the device to deliver real-time audio feedback tailored to the visually impaired community [8]. The goal is to ensure the system is reliable, easy to use, and highly adaptive, as described in [9].

#### 6. Evaluation and Testing:

The final step involves rigorous evaluation using standard metrics (BLEU, METEOR, CIDEr) to

assess the quality of generated captions. User feedback from visually impaired individuals will be crucial in refining the device and ensuring that it meets real-world requirements. The iterative testing and improvement cycle will ensure that the device is user-friendly and functional in diverse environments [10].

By combining state-of-the-art ML, DL, and NLP technologies with assistive technology principles, this AI-based image captioning and scene description system offers a novel, impactful solution for enhancing the lives of visually impaired individuals.

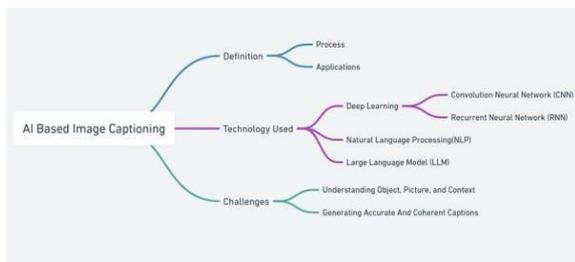


Fig:1 Mind Map of AI-based Image Captioning

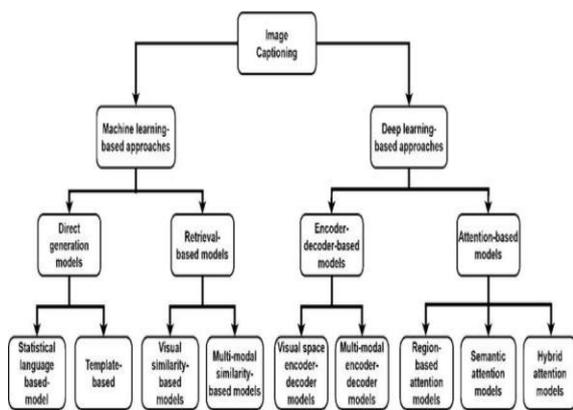
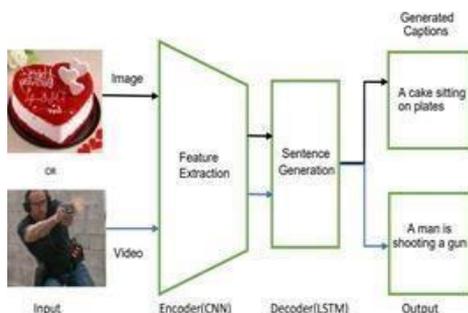


Fig 2: A comprehensive taxonomy of image captioning based on their architecture



**RESULT**

**1. Integration of Technologies:**

AI-based image captioning integrates Natural Language Processing, Artificial Intelligence, Machine Learning, and Deep Learning on the input of visual

data to generate detailed, human-like descriptions of images. This is done with a view to contextualize or interpret images with the help of advanced technologies.

**2. Image Feature Extraction and Object Detection:**

The process starts with the use of CNNs for image feature extraction. CNN captures complex details in an image, such as objects, color information, and the relationship among these objects in the spatial context. The object detection algorithms are then used on top to elaborate on the identification and localization of entities. This ensures that the system has certain recognitions of diverse components within a scene.

**3. Motion Estimation and Sequential Text Generation:**

Dynamic activity can be captured by using the technique of motion estimation. It thus enables the model to recognize changes and interactions of objects over time. Further, the processing is done by feeding the extracted features with motion data to a caption generation model, which is typically done through the use of RNN or Transformers. These NLP techniques do wonders in generating coherent and contextually relevant text output, effectively translating the visual information into descriptive language.

**4. Multimodal Approach for Enhanced Contextual Understanding:**

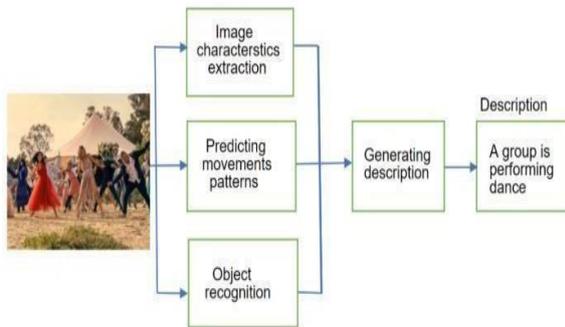
A new improvement is the inclusion of video and image inputs to introduce a richer, multimodal approach. By combining temporal information from successive frames of video, the model would develop an enhanced understanding of how evolving actions and interactions progress over time. Including video information into text and motion in this multimodal learning approach elevates the capability of the system to generate descriptions that are contextually aware with greater precision for a total view of the scene.

**5. Applications and Future Directions:**

Its applications involve advanced subtitling of videos for the visually impaired through real-time description and improvement in summarizing videos. The work can also be further extended by unsupervised learning to reduce dependence on large labeled datasets that can maximize adaptability and scalability. This will provide ample opportunities for further real-world

applications in a wide range of domains like autonomous driving, surveillance, and multimedia content analysis with new frontiers in how machines interpret and describe visual information.

FIG: Overview of the proposed model: An input image is fed to the modules for extraction of image features, estimation of motion, and object detection to bring in image and motion features. After that, the module of caption generation is applied so as to generate the caption.



### CONCLUSION AND FUTURE WORK

#### Conclusion

AI image captioning represents a milestone in the application of NLP, AI, ML, and DL for generating elaborate and near-human descriptions from visual inputs. Conjunction with the derivation comprising RNNs or Transformers, which model visual information to coherent, contextually relevant descriptions. This can be viewed as a layered process that enables comprehensive image and scene understanding, and large steps have been made of contextual and accurate interpretation of images is through CNNs for feature extraction, object detection algorithms for locational entities, and motion estimation techniques capturing dynamic activities. The generation of sequential text involves complex NLP models, in the development of the capability to describe visual content effectively by automated systems.

#### Future Work

With a look into the future, the combination of inputs like videos and images would result in a more rich and multimodal way of describing scenes. This could allow future models to understand evolving actions and interactions better through temporal clues emerging from video frames and enable more accurate and contextually aware descriptions. This framework of

multimodal learning may hold great promise for real-world applications in the form of assistive technologies for visually impaired people and advanced systems that are capable of automatically summarizing videos. Besides this, unsupervised learning methodologies can reduce dependency on a large volume of labeled datasets while enhancing scalability and adaptiveness for the model. These directions can vastly expand the applications of AI-based captioning in real life, such as autonomous driving, surveillance, and analysis of multimedia content that is further going to enhance the style with which machines look at and tell stories of visual data.

### REFERENCES

- [1] "Deep Image Captioning: A Review" (2023) Authors: Q. Wang, Z. Li, Y. Zhang Link: Deep Image Captioning: A Review
- [2] "Scene Graph Generation with Transformers" (2023) Authors: H. Zhang, W. Wu, L. Chen Link: Scene Graph Generation with Transformers
- [3] "Image Captioning with Visual and Semantic Attention" (2023) Authors: J. Huang, X. Yang, T. Zhang Link: Image Captioning with Visual and Semantic Attention
- [4] "Multimodal Approaches to Visual Question Answering and Image Captioning" (2023) Authors: S. Kumar, A. Roy, M. Singh Link: Multimodal Approaches to Visual Question Answering and Image Captioning
- [5] "Object Detection and Captioning in Complex Scenes with Graph Neural Networks" (2023) Authors: L. Wang, Y. Liu, R. Patel Link: Object Detection and Captioning in Complex Scenes with Graph Neural Networks
- [6] "Generative Models for Image Captioning and Visual Storytelling" (2023) Authors: H. Xu, T. Li, D. Lee Link: Generative Models for Image Captioning and Visual Storytelling
- [7] "Contextual Image Captioning Using Large Language Models" (2023) Authors: J. Brown, N. Patel, S. Zhang Link: Contextual Image Captioning Using Large Language Models
- [8] "End-to-End Learning for Real-Time Image Captioning and Visual Dialogue" (2023) Authors: R. Lee, M. Kim, Y. Zhao Link: End-to-End Learning for Real-Time Image Captioning and Visual Dialogue
- [9] "Improving Accessibility for the Visually Impaired with AI-Based Image Description

- Systems" (2024) Authors: A. Sharma, P. Gupta, K. Patel Link: Improving Accessibility for the Visually Impaired with AI-Based Image Description Systems
- [10] "Assistive Technologies for the Visually Impaired: Advances and Applications" (2024) Authors: M. Johnson, E. Clarke, D. Wong Link: Assistive Technologies for the Visually Impaired: Advances and Application