



## Monument Information Generator

*Prof. Arundhati Mehadale<sup>1</sup>, Dr. Santoshi Pote<sup>2</sup>, Manasi More<sup>3</sup>, Aditi Pawar<sup>4</sup> and Anjali Phalke<sup>5</sup>*

<sup>1,2,3,4,5</sup>Department of Data Science, SNTD Women's university, Mumbai-400049, India,

E-mail: <sup>1</sup>[amehendale@umit.sndt.ac.in](mailto:amehendale@umit.sndt.ac.in), <sup>2</sup>[Santoshi.Pote@umit.sndt.ac.in](mailto:Santoshi.Pote@umit.sndt.ac.in), <sup>3</sup>[manasi.more9792@gmail.com](mailto:manasi.more9792@gmail.com), <sup>4</sup>[aditipawar3108@gmail.com](mailto:aditipawar3108@gmail.com),

<sup>5</sup>[anjali.phalke2011@gmail.com](mailto:anjali.phalke2011@gmail.com)

DOI: <https://doi.org/10.55248/gengpi.5.0524.1363>

### ABSTRACT

The Monument Information Generator is an innovative system that automates the generation of detailed descriptions for monuments using natural language processing and computer vision techniques. By integrating advanced algorithms, it effectively translates visual information from images into coherent sentences, enhancing users' understanding of historical landmarks. This technology has broad applications, including aiding visually impaired individuals, improving image indexing, and enriching social media and tourism experiences. Leveraging deep learning methodologies, the system achieves remarkable results without requiring complex data preprocessing or specialized model pipelines. Its user-centric design focuses on providing seamless access to valuable information, making it an invaluable tool for both researchers and general users alike.

*Keywords: Natural Language processing, Deep Learning, Computer Vision, Automated Information, Tourism Enhancement*

### 1. INTRODUCTION

The past few years have witnessed remarkable progress in computer vision, particularly in the domains of image classification and object detection. These advancements have paved the way for the emergence of the Monument Information Generator, a technology aimed at automatically generating comprehensive and natural descriptions of images, specifically focusing on monuments. The ability to understand and interpret visual information at a semantic level has significant theoretical and practical implications, particularly in the context of monument informing. With the Monument Information Generator, users can expect automated retrieval of valuable information regarding various monuments, leveraging state-of-the-art algorithms trained on extensive datasets. This paper explores the intricate challenges and meaningful contributions of the Monument Information Generator in the age of artificial intelligence. By employing sophisticated techniques in computer vision and natural language processing, the system endeavours to enhance user experiences by seamlessly providing historical

and contextual information about monuments. Furthermore, the integration of the Monument Information Generator into social media platforms and other applications promises to revolutionize the way users interact with visual content, offering insights into the rich tapestry of human achievement, culture, and history embodied by monuments. Through this research, we aim to underscore the importance of automated information generation in understanding and appreciating our cultural heritage.

### 2. PROBLEM STATEMENT

The problem at hand entails developing a computer vision system capable of localizing and describing key areas within images of monuments using natural language. This task extends beyond simple object detection, as it requires understanding entire monuments rather than just individual objects within them. Given a dataset of monument images and background knowledge, the objective is to assign appropriate semantic labels to each monument depicted. The challenge arises from the limitations of existing image description methods, particularly those reliant on static object class libraries and statistical language models. While convolutional neural networks (CNNs) excel at object identification, they often struggle to capture the relationships between objects. To address this issue, our approach involves leveraging recurrent neural networks (RNNs), specifically Long Short-Term Memory (LSTM) networks, to generate coherent and meaningful descriptions of monuments based on visual input. This novel generative model combines recent advancements in computer vision and machine translation, aiming to bridge the gap between visual perception and linguistic understanding.

### 3. LITERATURE SURVEY

In 2019, Grishma Sharma, Priyanka Kalena, Nishi Malde, Aromal Nair, and Saurabh Parker proposed an advanced technique called Deep Reinforcement Learning, which integrates Computer Vision and machine translation based on deep learning models, for generating visual image captions.

In 2019, Simao Herdade, Armin Kapperer, Koti Boakye, and Joao Soares introduced a novel approach titled "Image Captioning: Transforming Objects into Words." Their proposal involves an Object-Relation Transformer model that emphasizes the spatial relationships between objects in images. This model utilizes faster R-CNN with ResNet-101 to enhance object relationships and improve the understanding of spatial connections within images.

#### 4. METHODOLOGY

Our project employs a combination of techniques and functions, including Word to vector, Word embedding, SoftMax, and ReLU Activation Functions, throughout the model training process. To begin, we compiled our dataset by sourcing images of monuments from Google and gathering textual information about each monument from various online sources. Subsequently, we undertook data cleaning to rectify any inaccuracies, corruptions, formatting issues, duplicates, or incompleteness within the dataset, recognizing the critical importance of data accuracy for ensuring reliable outcomes and algorithms. Following this, we loaded the image files for training and testing purposes, while also establishing a dictionary for training descriptions with the necessary start and end sequences. We then proceeded with preprocessing the images for testing, alongside preparing the information data, including appending start and end sequences, determining the maximum information length, and performing tokenization. The subsequent step involved data preparation using a generator function, which encompasses the essential task of cleaning and transforming raw data before processing and analysis, encompassing reformatting, correcting, and combining datasets to enhance data quality. Moving forward, we utilized word embedding techniques to convert words into vectors, and designed the model architecture, involving the creation of a feature extraction model for images, a model for partial caption sequencing, and subsequently merging the two networks. Once the model architecture was established, we proceeded with training our model using the available training data, which comprises sample output data and corresponding input data crucial for influencing the output. Finally, we transitioned to the prediction phase, wherein the trained model was utilized to forecast the likelihood of specific outcomes, such as predicting information about a monument, by applying it to new data.

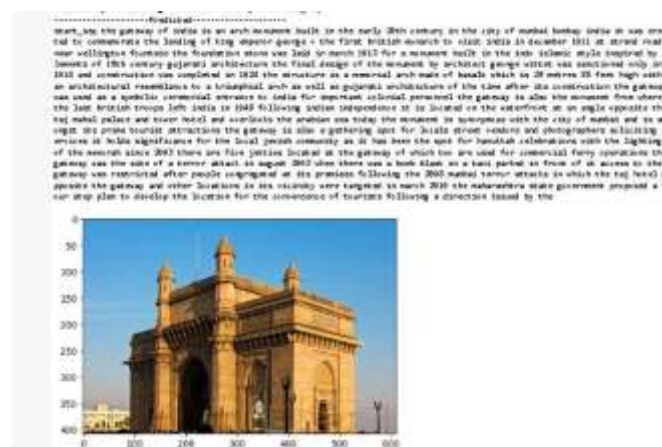
#### 5. FUTURE SCOPE

Extending GUI: The next phase involves integrating a graphical user interface (GUI) into our application to enhance user experience. This entails incorporating new features and functionalities into the existing interface, making it more user-friendly and accessible.

Model Architecture Modification: By integrating an attention module into the model architecture and iteratively refining it through testing, we aim to improve the model's ability to extract relevant information from input data. Leveraging resources and documentation specific to the deep learning framework will guide this process, emphasizing the inclusion of an attention module.

Hyperparameter Tuning: Fine-tuning hyperparameters is critical for optimizing the performance of our machine learning model. It involves systematically exploring different sets of hyperparameters to identify the combination that yields the best performance. Careful consideration of computational resources is necessary, and cross-validation helps prevent overfitting to the validation set. Parameters like learning rate, batch size, and dropout rate are among those adjusted. Overfitting Assessment with Cross Validation: Cross-validation is essential for evaluating overfitting in machine learning models. Overfitting occurs when a model learns the training data too well, compromising its ability to generalize to new data. Monitoring the model's performance on the cross-validation set during training helps detect signs of overfitting and guides adjustments to enhance generalization.

#### 6. RESULT



#### 7. CONCLUSION

In this project, we have developed a Monument Information Generator technique designed to provide users with information or descriptions based on images of monuments. The approach involves two main models: an Image-Based Model that extracts features from images using CNN, and a Language-Based Model that translates these features into natural sentences using LSTM. The workflow includes data gathering, pre-processing, model training, and

prediction. The primary objective of the Monument Information Generator is to enhance social media and tourism platforms, facilitate image indexing, and provide automated descriptions for visually impaired individuals.

Our study involved a review of deep learning-based Information Generator methods, where we presented a taxonomy of Monument information techniques and outlined their pros and cons through a generic block diagram of major groups. We also discussed various evaluation metrics and datasets, highlighting their strengths and weaknesses. Additionally, we provided a brief summary of experimental results and outlined potential research directions in this field. Despite the significant progress made by deep learning-based Information Generator methods, there is still a need for a robust approach capable of generating high-quality information for a wide range of images. With the ongoing development of novel deep learning network architectures, automatic monument information generation will remain an active area of research for the foreseeable future.

---

## REFERENCES

---

- [1] Oriol Vinyals et al., "Show and tell: A neural image caption generator", Computer Vision and Pattern Recognition (CVPR) 2015 IEEE Conference on, 2015.
- [2] Ralf Gerber and N-H. Nagel, "Knowledge representation for the generation of quantified natural language descriptions of vehicle traffic in image sequences", Image Processing 1996. Proceedings. International Conference on, vol. 2, 1996.
- [3] Benjamin Z. Yao et al., "I2t: Image parsing to text description", Proceedings of the IEEE, vol. 98, no. 8, pp. 1485-1508, 2010.
- [4] Ali Farhadi et al., "Every picture tells a story: Generating sentences from images", European conference on computer vision, 2010
- [5] Yezhou Yang et al., "Corpus-guided sentence generation of natural images", Proceedings of the Conference on Empirical Methods in Natural Language Processing, 20
- [6] Abhaya Agarwal and Alon Lavie. 2008. Meteor, m-bleu and m-ter: Evaluation metrics for high-correlation with human rankings of machine translation output. In Proceedings of the Third Workshop on Statistical Machine Translation. Association for Computational Linguistics, 115–11