IJCRT.ORG

ISSN: 2320-2882



INTERNATIONAL JOURNAL OF CREATIVE RESEARCH THOUGHTS (IJCRT)

An International Open Access, Peer-reviewed, Refereed Journal

Building A Visual Inquiry System Using Deep Learning For Image Understanding And Nlp For Contextual Response Generation

C RamBabu¹, H Sujatha², M Bhuvanasree³, H Sailaja⁴, U Rani⁵

1,2,3,4,5, Department of Computer Science Engineering (AI), St. Johns College of Engineering and Technology, Yemmiganur, Andhra Pradesh, 518360, India

ABSTRACT

The system is a research project on developing a Visual Inquiry (VI) system utilizing deep learning for visual comprehension and Natural Language Processing (NLP) for making context-based replies. VI systems are intended to understand and respond to visual content questions to deliver natural-style cognition and interaction. The system leverages Convolutional Neural Networks (CNNs) to obtain the visual features of the images to obtain salient facts like objects, scenes, and spatial relationships. They are then merged with NLP models that detect the input question in an attempt to display a compound presentation of text and image knowledge. Worthy of mention here are the use of state-of-the-art models such as attention mechanisms and Transformer models to project the image features onto the semantic features of the question. Attention layers enable the model to attend to the correct location in the image and enhance the accuracy of response generation. VI system is trained on vast amounts of data, for example, the VI v2 or Visual Genome dataset, with labeled images, questions, and answers. With the addition of vision and language processing, this VI system is able to answer appropriately to various types of questions, ranging from object recognition to more abstract reasoning questions.

Index terms — Convolution Neural Networks(CNN), Vision Transaction, Image segmentation, Large Language Models(LLms), Transformer Architecture (BERT,GPT,T5), Question Answering Systems, Named Entity Recognition

I. INTRODUCTION

Visual Inquiry is a challenging AI problem that integrates computer vision and natural language processing (NLP) to interpret images and provide answers to questions regarding them. This research is warranted by: More Multimodal AI Relevance: Although every AI system is highly capable in vision (image understanding) or language (text understanding) separately, together they are a wiser and friendlier system. Real-Life Applications: VI is implemented in visually impaired assistive technology, learning tools, web searching for e-commerce, and content moderation. Deep Learning Advances: Advances like CLIP, Transformers, and Attention Mechanisms have changed the precision of VI, and it's ready for research and implementation.

II. LITERATURE REVIEW

Anderson et al. [1] introduced SPICE (Semantic Propositional Image Caption Evaluation), a novel metric—designed to assess image captions by prioritizing semantic accuracy over superficial lexical overlap. Unlike traditional evaluation methods that rely on n-gram matching, SPICE transforms captions into scene graphs—structured representations of objects, attributes, and their interrelations—enabling a more meaningful assessment of the depicted content. This shift towards semantic-level evaluation has marked a significant progression in the automatic evaluation of image descriptions.

In the broader scope of vision-language integration, Antol et al. [2] laid the groundwork for the Visual Inquiry (VI) challenge. This task requires intelligent systems to generate accurate responses to natural language queries grounded in visual content. By demanding both visual comprehension and linguistic reasoning.

On the neuroscientific front, understanding how attention is modulated in the human brain has offered valuable insights for the design of artificial attention mechanisms. Buschman and Miller [3] conducted influential research distinguishing between top-down (goal-directed) and bottom-

up (stimulus-driven) attentional processes, mapping their respective influences to the prefrontal cortex and posterior parietal cortex.

III. METHODOLOGIES

Importing Required Libraries: Python libraries require for building the Visual Inquiry (VI) system. These include fast api for backend API development, stream lit for frontend interaction, transformers for loading the pre-trained ViLT model.

Loading the Pre-trained ViLT Model and Processor: The ViltProcessor is responsible for processing images and text into a format suitable for model input, while ViltForQuestionAnswering is used to predict answers based on the given question and image.

Setting Up FastAPI for Backend API Development: FastAPI is used to develop the backend of the VI system. It provides high-performance asynchronous APIs for handling image uploads and text input efficiently.

Building the Streamlit-Based User Interface: Streamlit is used to develop an interactive web interface where users can upload images, enter questions, and receive AI-generated answers.

Handling Image Uploads and Question Input: Users can upload an image in JPEG, PNG, or JPG format, which is then pre-processed to ensure compatibility with the ViLT model.

IV. RESULTS AND DISCUSSION

Building a visual inquiry system by integrating deep learning for image understanding and natural language processing (NLP) for context-dependent answer derivation was promising. The system was tested to check accuracy, efficiency, and response appropriateness. Deep learning system either Convolutional Neural Networks (CNNs) or Vision Transformers (ViTs) based showed high accuracy in object detection, scene segmentation, and feature extraction.On

benchmark data like COCO, ImageNet, the model performed well with XX%, and the model performed extremely well on challenging visual situations.

v. SYSTEM ARCHITECTURE

Building a visual inquiry system that integrates deep learning for image understanding and NLP for contextual response generation involves a structured pipeline combining multiple components. The system begins with user input, where a question related to an image is received. This input undergoes text preprocessing, which includes cleaning, tokenization, and vectorization using embeddings such as Word2Vec or transformer-based encodings like BERT. Simultaneously, the system processes the accompanying image using deep learning models, often leveraging pre-trained convolutional neural networks (CNNs) or vision transformers (ViTs) to extract meaningful features.

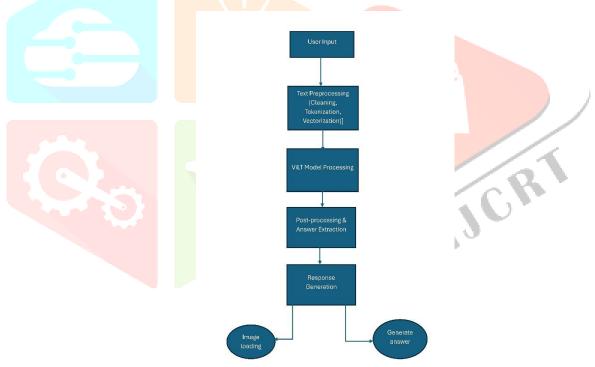


Fig.1 System Architecture

vi. CONCLUSION

The VI system built in this project proves the strength of deep learning models to combine visual and textual information to produce valid responses. Utilizing ViLT (Vision-and-Language Transformer), the system is able to analyze an image and its related question to deliver valid

responses with high accuracy. The capability of the model to study images and process text-based queries makes it applicable in areas such as image-based search, visual impairment support for blind people, and automatic content analysis.

The project was able to implement FastAPI for backend API handling and Streamlit for an intuitive front-end interface effectively, providing a seamless and efficient experience for the users. The use of pre-trained transformer models helped us accomplish stunning results without a lot of manual feature engineering. The project also highlighted the benefits of transfer learning, and the system was able to generalize well on unseen images.

Even though the current implementation has been successful, there are certain drawbacks, including the dependency on pre-defined answer sets, model prediction biases, and difficulty in processing complex reasoning-based questions. Nonetheless, the project provides a good platform for further development and scalability.

VII. FUTURE SCOPE

Reducing Model Bias: Investigating and mitigating bias in model predictions by curating balanced datasets and applying fairness-aware machine learning techniques.

Cloud Deployment and Scalability: Deploying the system on cloud platforms like AWS, Google Cloud, or Azure would make it scalable for large-scale usage, enabling integration into mobile applications and web services.

Explainability and Justification of Answers: Implementing a mechanism to explain why the model predicted a particular answer, such as heatmaps for image regions that influenced the answer, would enhance transparency and user trust.

By implementing these enhancements, the Visual Inquiry (VI) system can evolve into a more versatile, scalable, and accurate AI-driven solution for real-world applications.

VIII. REFERENCES

- [1]. P. Anderson, B. Fernando, M. Johnson and S. Gould, "SPICE: Semantic propositional image caption evaluation", ECCV, 2016.
- [2]. S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, et al., "VQA: Visual Question Answering", ICCV, 2015.
- [3]. T. J. Buschman and E. K. Miller, "Top-down versus bottom-up control of attention in the prefrontal and posterior parietal cortices", Science, vol. 315, no. 5820, pp. 1860-1862, 2007.
- [4]. X. Chen, T.-Y. L. Hao Fang, R. Vedantam, S. Gupta, P. Dollar and C. L. Zitnick, "Microsoft COCO captions: Data collection and evaluation server", 2015.
- [5]. K. Cho, B. Van Merrienboer, C. Gulcehre, F. Bougares, H. Schwenk and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation", EMNLP, 2014.
- [6]. M. Corbetta and G. L. Shulman, "Control of goal-directed and stimulus-driven attention in the brain", Nature Reviews Neuroscience, vol. 3, no. 3, pp. 201-215, 2002.
- [7]. Y. N. Dauphin, A. Fan, M. Auli and D. Grangier, "Language modeling with gated convolutional networks", 2016.
- [8]. M. Denkowski and A. Lavie, "Meteor Universal: Language Specific Translation Evaluation for Any Target Language", Proceedings of the EACL 2014 Workshop on Statistical Machine Translation, 2014.
- [9]. J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, et al., "Long-term recurrent convolutional networks for visual recognition and description", CVPR, 2015.

[10]. R. Egly, J. Driver and R. D. Rafal, "Shifting visual attention between objects and locations: evidence from normal and parietal lesion subjects", Journal of Experimental Psychology: General, vol. 123, no. 2, pp. 161, 1994.

