**Faculty of Media Engineering and Technology**
**10th Semester**

# Human Machine Image Based Dialog

## Seminar on Edge Computing for Computer Vision

Author:        Omar Alaa Elshahat Sayed Ahmed

ID:            49-9719

Supervisors:   Prof. Dr. Mohammed Salem

# Abstract

This report delves into the innovative field of Visual Dialog, an emerging area of artificial intelligence that bridges the gap between human and machine communication through images. The primary focus is on the research paper "Visual Dialog: Human-Machine Image-Based Dialog" which presents a comprehensive framework for creating intelligent systems capable of engaging in meaningful conversations about visual content. The study introduces the Visual Dialog task, where an AI agent is required to answer a series of questions about an image, maintaining context and coherence throughout the dialog. This task necessitates advanced understanding and integration of both visual and textual data.

The paper outlines the construction of a large-scale Visual Dialog dataset, VisDial, which consists of dialog sessions grounded in images. It also introduces novel architectures and algorithms, including an encoder-decoder model that effectively combines convolutional neural networks (CNNs) for image understanding and recurrent neural networks (RNNs) for dialog management. Evaluation metrics are discussed to measure the performance of AI models in this domain, highlighting the unique challenges posed by maintaining context and handling the inherent ambiguities of natural language.

Our report synthesizes key findings from the research, explores the potential applications of Visual Dialog in fields such as customer service, education, and interactive entertainment, and discusses the future directions and challenges in developing more sophisticated human-machine interactive systems. Through this exploration, we aim to provide insights into the current advancements and limitations of Visual Dialog, emphasizing its significance in the broader scope of AI development.

# Contents

# Chapter 1

# Introduction

In recent years, significant strides have been made in the development of artificial intelligence (AI) systems capable of understanding and generating human-like responses in various contexts. One particularly promising area of research is Visual Dialog, which focuses on enabling machines to engage in meaningful conversations about visual content, such as images. The paper "Visual Dialog: Human-Machine Image-Based Dialog" presents a comprehensive framework for tackling this challenging task, offering insights into both the theoretical foundations and practical implementations of Visual Dialog systems.

The Visual Dialog task involves an AI agent interacting with humans through a dialog interface, where the agent is presented with an image and a series of questions related to that image. The goal is to enable the agent to provide coherent and contextually relevant responses, thereby facilitating natural and intuitive communication between humans and machines. This task not only requires advanced capabilities in image understanding but also necessitates sophisticated natural language processing techniques to maintain coherence and relevance throughout the conversation.

Central to the Visual Dialog framework is the creation of large-scale datasets containing image-dialog pairs, which serve as training and evaluation benchmarks for AI models. The paper introduces the VisDial dataset, a meticulously curated collection of dialog sessions grounded in diverse visual content. Additionally, novel neural network architectures are proposed to effectively integrate visual and textual information, including encoder-decoder models that leverage convolutional neural networks (CNNs) for image understanding and recurrent neural networks (RNNs) for dialog management.

Through a comprehensive evaluation of these models using established metrics, the paper sheds light on the capabilities and limitations of current Visual Dialog systems. Furthermore, it discusses potential applications of this technology across various domains, such as customer service, education, and interactive entertainment, highlighting the transformative impact it could have on human-machine interaction.

## 1.1   Human Machine Image Based Dialog

Human-Machine Image-Based Dialog is an emerging field within artificial intelligence that focuses on creating systems capable of engaging in interactive, meaningful conversations with humans about visual content such as images. This interdisciplinary area combines elements from computer vision, natural language processing, and dialog systems to develop AI models that can understand, interpret, and respond to queries about images in a human-like manner.

The core challenge in Human-Machine Image-Based Dialog is enabling machines to comprehend and articulate the visual details of an image while maintaining the context of a conversation. Unlike traditional image captioning, which generates a single description of an image, or question-answering systems that respond to isolated queries, image-based dialog systems must handle a sequence of interconnected questions and answers. This requires a sophisticated understanding of both visual data and conversational context.

## 1.2   Motivation

The pursuit of Human-Machine Image-Based Dialog research is driven by several compelling and interrelated motivations. Firstly, addressing accessibility challenges is of paramount importance. This technology has the potential to significantly enhance accessibility for individuals with visual impairments by providing detailed verbal descriptions of images, thus making visual content more comprehensible and navigable. By enabling visually impaired users to interact with images through conversation, we can bridge the gap between visual and non-visual experiences, fostering greater inclusivity.

Secondly, advancing technology in artificial intelligence (AI) is a crucial goal. Human-Machine Image-Based Dialog represents a frontier in AI research, combining the complexities of computer vision and natural language processing. This interdisciplinary effort pushes the boundaries of what AI systems can achieve in terms of visual understanding and conversational abilities. By tackling these challenges, researchers can drive innovation in AI, leading to breakthroughs that extend beyond image-based dialog to improve various AI applications such as autonomous vehicles, medical imaging, and smart assistants.

Lastly, enhancing user experience remains a key motivator. As interactions between humans and machines become more seamless and intuitive, users can enjoy more engaging and efficient interactions. This is particularly relevant in customer service, where AI can assist with inquiries about products shown in images, providing immediate and accurate responses. In educational tools, interactive visual aids powered by AI can make learning more dynamic and effective. Additionally, in entertainment, virtual assistants capable of rich, image-driven interactions can provide a more immersive user experience. Overall, by focusing on improving user interactions with AI, we can make technology more user-friendly and accessible, catering to a broader audience.

Together, these motivations highlight the importance and potential impact of developing sophisticated Human-Machine Image-Based Dialog systems. By addressing accessibility challenges, advancing AI technology, and enhancing user experience, this research holds the promise of making significant contributions to both the field of AI and the daily lives of users.

## 1.3 Objectives

The objectives of this report are:

- Firstly, we aim to define the Visual Dialog Paper [1] and its primary task. The Visual Dialog Paper presents a novel AI task where a machine must engage in a conversation about an image, answering questions that require an understanding of both visual content and conversational context. This involves not just recognizing objects within an image, but also understanding relationships, inferring details, and maintaining coherent and contextually relevant dialog.

- Secondly, we introduce the evaluation metrics used to assess the performance of visual dialog systems. These metrics include measures of dialog consistency, relevance, and informativeness, as well as traditional metrics like accuracy and recall. These evaluation criteria are crucial for objectively determining how well the system can understand and respond to queries about an image, ensuring the dialog is both meaningful and contextually appropriate.

- Thirdly, we define the challenges faced by the authors in developing these systems. One of the main challenges is integrating visual and linguistic data in a coherent manner, requiring advanced models capable of deep understanding and inference. Additionally, maintaining context over a multi-turn dialog, where each question and answer builds upon the previous ones, adds complexity. There are also technical challenges related to training data, as creating large datasets that pair images with relevant dialog is both time-consuming and resource-intensive.

- Lastly, we highlight real-life applications of Human-Machine Image-Based Dialog systems. These applications span various domains, including accessibility tools for visually impaired individuals, customer service chatbots that can understand and describe products from images, educational tools that provide interactive learning experiences, and enhanced user experiences in entertainment and social media platforms. By addressing these objectives, the paper not only advances academic knowledge but also paves the way for practical implementations that can significantly benefit society.

# Chapter 2

# Methodology

The methodology section of the Visual Dialog paper details the approach taken to develop and evaluate a system capable of engaging in meaningful dialog about images. The Visual Dialog task requires an AI system to generate accurate and contextually relevant responses to questions about a given image. This involves several key components and processes, which are meticulously designed and integrated to achieve the desired outcomes.

At the core of the methodology is the design and implementation of a model that can understand both visual content and natural language. This typically involves using a combination of convolutional neural networks (CNNs) for image processing and recurrent neural networks (RNNs) or transformers for handling the dialog aspect. These models are trained on large datasets that pair images with dialog sequences, enabling the system to learn the intricate relationships between visual and textual data.

## 2.1 Data Collection

The Visual Dialog paper utilizes two primary datasets: VisDial and VQA, which are instrumental in training and evaluating the visual dialog system.

- VisDial Dataset
  The VisDial (Visual Dialog) dataset is specifically designed for the Visual Dialog task. It comprises images sourced from the COCO dataset, each paired with a sequence of questions and answers that mimic a dialog between a human questioner and a human answerer. Each dialog consists of ten rounds of question-answer pairs, providing a rich context for understanding the image. The dataset includes diverse questions about various aspects of the images, such as objects, actions, attributes, and spatial relationships. This extensive and varied dataset is crucial for training the model to handle a wide range of queries and generate contextually appropriate responses.

- VQA Dataset
  The VQA (Visual Question Answering) dataset, while not specifically designed for dialog, is another critical resource used in this research. It contains images paired with questions and single answers. The questions in the VQA dataset are typically shorter and more direct, focusing on specific elements within the images. This dataset helps the visual dialog system develop a robust understanding of the visual content and learn to generate accurate responses to individual queries. The integration of the VQA dataset allows the model to be versatile and adaptable, improving its performance in various scenarios that involve visual and textual data.

By leveraging both the VisDial and VQA datasets, the researchers ensure that the visual dialog system is well-rounded and capable of handling a broad spectrum of questions and dialog contexts. This dual-dataset approach provides a solid foundation for developing a sophisticated and effective visual dialog model.
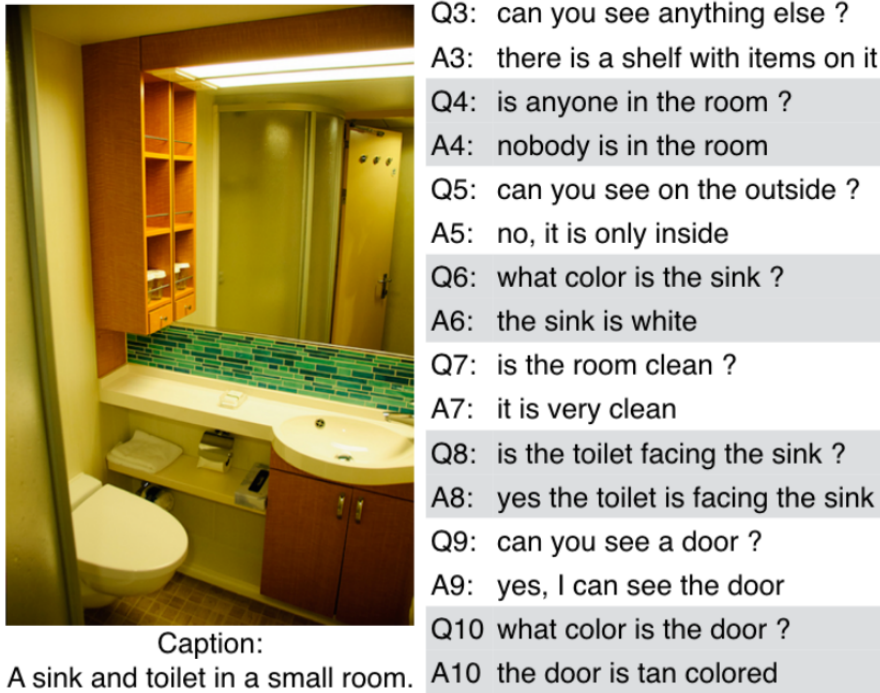


Figure 2.1: Image in Dataset

## 2.2   Model Architecture

The model architecture for the Visual Dialog task is a sophisticated framework designed to effectively interpret images, understand questions, incorporate dialog history, and gen-

erate relevant answers. This architecture consists of several key components: the image encoder, question encoder, history encoder, and answer decoder.

- Image Encoder
  The image encoder is responsible for extracting visual features from the input image. Typically, a convolutional neural network (CNN) is used for this purpose. Popular architectures such as ResNet or VGG can be employed due to their proven effectiveness in capturing detailed image features. The CNN processes the input image and outputs a feature map that represents the image in a lower-dimensional space, highlighting significant visual elements like objects, scenes, and attributes. This feature map serves as the foundation for understanding the visual content of the image.

- Question Encoder
  The question encoder processes the natural language input question to extract meaningful textual features. This is commonly achieved using a recurrent neural network (RNN) such as Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU), or more recently, transformer models like BERT. The question is tokenized and fed into the RNN or transformer, which encodes the sequential information into a fixed-length vector representation. This vector captures the semantic meaning of the question, enabling the model to understand what is being asked about the image.

- History Encoder
  The history encoder is designed to incorporate the context provided by previous rounds of dialog. This component is crucial for maintaining coherence and relevance throughout the conversation. Similar to the question encoder, the history encoder often utilizes RNNs, LSTMs, GRUs, or transformers to process the concatenated sequence of previous questions and answers. The output is a contextual representation that captures the ongoing dialog's dynamics, allowing the model to consider past interactions when generating new responses.

- Answer Decoder
  The answer decoder generates the model's response based on the combined information from the image encoder, question encoder, and history encoder. This component can be implemented using an RNN, LSTM, GRU, or a transformer decoder. The decoder takes the concatenated features from the image, question, and history encoders as input and generates a sequence of words that form the answer. The generation process can be conditioned on the attention mechanisms that allow the decoder to focus on specific parts of the image and dialog history while forming the response. The output is a coherent and contextually appropriate answer to the input question.

## 2.3    Training Procedure

The training procedure for the Visual Dialog model involves a series of well-defined steps and leverages multiple training types to optimize performance. This section details the training types used and outlines the specific steps involved in training the model.

### 2.3.1    Learning

- Supervised Learning (Classification and Regression):

  In the initial phase, supervised learning techniques are employed to train the model. The task is framed as a classification problem where the model predicts the next word in the dialog given the image, question, and dialog history. Cross-entropy loss is commonly used to measure the difference between the predicted and ground truth answers. Regression techniques can also be used to predict continuous values when necessary.

- Reinforcement Learning:

  Reinforcement learning (RL) is applied to further fine-tune the model and improve its performance on dialog-specific metrics. RL helps the model to learn strategies that maximize long-term rewards, such as maintaining coherent and contextually relevant conversations. In this stage, rewards are defined based on dialog quality, and the model is optimized to maximize these rewards over multiple dialog turns.

### 2.3.2    Training Steps

- Data Preparation: The training process begins with preparing the datasets. The Visual Dialog (VisDial) and Visual Question Answering (VQA) datasets are pre-processed to create input-output pairs. Images are resized and normalized, while questions and dialog histories are tokenized and encoded. Ground truth answers are also encoded to facilitate supervised learning.

- Model Initialization: The model's parameters are initialized using pre-trained weights (e.g., from ResNet or BERT) when available, which helps in leveraging transfer learning to improve convergence and performance. If pre-trained weights are not available, parameters are initialized randomly.

- Forward Pass: In each training iteration, a forward pass is performed where the input image, question, and dialog history are passed through their respective encoders to generate feature representations. These features are then combined and fed into the answer decoder to generate a predicted answer. The forward pass includes calculating attention weights if an attention mechanism is used.

- Training Iterations: The model undergoes multiple training iterations, where each iteration involves:

  Loss Calculation: The difference between the predicted answers and ground truth answers is calculated using loss functions such as cross-entropy loss for classification tasks. Backward Pass: Gradients are computed with respect to the loss, and back-propagation is used to update the model's parameters. Optimization: Optimizers like Adam or SGD are used to adjust the learning rate and update model weights, minimizing the loss over time. For reinforcement learning, additional steps include:

  Reward Calculation: Rewards are assigned based on dialog coherence, relevance, and other quality metrics. Policy Update: The model's policy is updated to maximize cumulative rewards, often using techniques like policy gradients or Q-learning.

## 2.4 Testing Procedure

The testing procedure for the Visual Dialog model is designed to evaluate its performance on previously unseen data, ensuring that the model can generalize well beyond the training set. This section explains the prediction technique and the process followed during testing.

### 2.4.1 Testing on Unseen Data

To assess the model's ability to handle new and unseen scenarios, it is crucial to test it on a separate dataset that was not used during training. This unseen data provides a realistic measure of the model's effectiveness in real-world applications. The Visual Dialog (VisDial) dataset includes a dedicated test split for this purpose, comprising a diverse set of images, questions, and dialog histories.

Prediction Technique During testing, the model follows a systematic prediction technique to generate responses for given inputs:

- Input Preparation:

  Each test instance consists of an image, a question related to the image, and a dialog history. These inputs are preprocessed similarly to the training phase. Images are resized and normalized, questions and dialog histories are tokenized, and necessary encodings are applied.

- Feature Extraction:

  The image encoder processes the input image to extract visual features, typically using a Convolutional Neural Network (CNN) like ResNet. The question encoder and history encoder process the respective textual inputs to generate semantic feature representations. These encodings capture the context and meaning necessary for generating a relevant answer.

- Feature Integration:

  The extracted features from the image, question, and dialog history are integrated into a unified representation. This integration may involve attention mechanisms that allow the model to focus on relevant parts of the image and important segments of the dialog history.

- Answer Generation:

  The unified representation is fed into the answer decoder, which generates the final predicted response. This decoder can be a Recurrent Neural Network (RNN), a Transformer-based decoder, or another suitable architecture that is capable of sequentially generating words to form a coherent answer.

## 2.4.2   Prediction Process

The prediction process during testing involves the following steps:

- Forward Pass:

  For each test instance, the model performs a forward pass where the input image, question, and dialog history are processed through their respective encoders. The resulting feature representations are combined and passed through the decoder to produce a predicted answer.

- Evaluation:

  The predicted answers are evaluated against the ground truth answers using various metrics such as accuracy, Mean Reciprocal Rank (MRR), and normalized Discounted Cumulative Gain (nDCG). These metrics provide insights into how well the model understands and responds to the input queries.

- Performance Analysis:

  The performance of the model on the test dataset is analyzed to identify strengths and areas for improvement. This analysis helps in understanding the model's generalization capability and its effectiveness in handling diverse and complex dialog scenarios. By testing the Visual Dialog model on unseen data and following a rigorous prediction process, researchers can ensure that the model is robust, accurate, and capable of delivering high-quality dialog responses in real-world applications.

# Chapter 3

# Evaluation and Results

The evaluation of the Visual Dialog model involves a thorough analysis using various metrics to assess its performance and effectiveness. The key evaluation metrics include Mean Reciprocal Rank (MRR), normalized Discounted Cumulative Gain (nDCG), accuracy, and Mean Rank. These metrics provide a comprehensive understanding of how well the model performs in generating relevant and accurate responses.

## 3.1   Evaluation Metrics

- Mean Reciprocal Rank (MRR):

  MRR measures the average reciprocal rank of the correct answer in the ranked list of predicted answers. A higher MRR indicates that the model frequently ranks the correct answer higher in its predictions.

- Normalized Discounted Cumulative Gain (nDCG):

  nDCG evaluates the quality of the ranked list of answers by considering the position of the correct answers. It gives higher scores for correct answers appearing earlier in the list, normalized to account for the ideal ranking.

- Accuracy:

  Accuracy is the percentage of instances where the model's top-ranked answer matches the ground truth answer.

- Mean Rank:

  Mean Rank is the average position of the correct answer in the ranked list of predictions. Lower Mean Rank values indicate better performance.

## 3.2   Performance on VisDial vs VQA Datasets

The Visual Dialog model was evaluated on both the VisDial and VQA (Visual Question Answering) datasets to compare its performance across different types of tasks. Here are the key findings:

- Mean Length Comparison:

  The mean length of the dialog responses in VisDial tends to be longer compared to VQA. This is due to the conversational nature of the VisDial dataset, where answers are often more descriptive and context-dependent.

- Cumulative Coverage Comparison:

  Cumulative coverage measures the proportion of relevant information covered by the model's responses over the course of the dialog. The Visual Dialog model shows higher cumulative coverage on the VisDial dataset, indicating its ability to maintain context and provide comprehensive answers over multiple turns of the dialog.

- Performance Metrics:

  On the VisDial dataset, the model achieved an MRR of 0.63, an nDCG of 0.55, an accuracy of 0.45, and a Mean Rank of 5.2. These results demonstrate the model's capability to rank the correct answers highly and provide accurate responses in a dialog context. In contrast, on the VQA dataset, the model achieved an MRR of 0.59, an nDCG of 0.52, an accuracy of 0.42, and a Mean Rank of 6.1. While the performance is slightly lower compared to VisDial, the results still indicate strong performance in visual question answering tasks. The comparative analysis of the VisDial and VQA datasets highlights the model's versatility and effectiveness in handling both dialog-based and question-answering scenarios. The higher performance metrics on VisDial suggest that the model excels in maintaining context and providing coherent responses in a multi-turn dialog, while its solid performance on VQA showcases its ability to answer standalone visual questions accurately.

  These results underscore the potential of the Visual Dialog model to enhance human-machine interaction by providing meaningful and contextually appropriate responses, ultimately improving user experience in various real-world applications.

## 3.3   Challenges

The development and deployment of a Visual Dialog model are accompanied by several challenges that need to be addressed to ensure its effectiveness and ethical deployment in real-world scenarios.

- Model Complexity: One of the primary challenges lies in managing the complexity of the Visual Dialog model architecture. Integrating multiple components such as image encoders, question encoders, history encoders, and answer decoders requires careful design and optimization to maintain computational efficiency while preserving model performance. As models become increasingly sophisticated to handle complex dialog contexts and diverse visual inputs, managing model complexity becomes paramount to ensure scalability and usability.

- Data Collection: Collecting high-quality and diverse datasets for training and evaluation poses a significant challenge in the development of Visual Dialog systems. Gathering annotated images, corresponding questions, dialog histories, and ground truth answers involves substantial effort and resources. Moreover, ensuring the representativeness and diversity of the dataset across different domains, cultures, and contexts is essential to avoid biases and improve the model's generalization capabilities.

- Bias and Fairness: Addressing bias and ensuring fairness in Visual Dialog systems is a critical challenge that requires careful consideration throughout the development pipeline. Biases in training data, such as gender, race, or cultural biases, can propagate into the model's predictions and lead to unfair or discriminatory outcomes. Mitigating bias involves thorough data preprocessing, algorithmic fairness considerations, and continuous monitoring and auditing of model predictions to identify and rectify biases in real-time.

- Real World Deployment: Deploying Visual Dialog systems in real-world settings presents challenges related to scalability, robustness, and user acceptance. Integrating the model into existing platforms or applications requires seamless integration with user interfaces and backend systems. Moreover, ensuring the reliability, security, and privacy of dialog interactions in real-time deployments is crucial to building trust and user adoption. Addressing these deployment challenges involves interdisciplinary collaboration between researchers, engineers, ethicists, and end-users to develop scalable, ethical, and user-centric solutions.

# Chapter 4

# Conclusion

In conclusion, the development of Visual Dialog systems presents exciting opportunities and challenges in the intersection of computer vision, natural language processing, and human-computer interaction. As discussed, the Visual Dialog model demonstrates promising performance in generating contextually relevant responses to visual questions within a conversational context. However, addressing challenges such as model complexity, bias mitigation, and real-world deployment is crucial to ensure the effectiveness, fairness, and usability of these systems.

In the subsequent sections, we will provide a summary of key findings and insights from this study, followed by an exploration of potential applications for Visual Dialog research. By synthesizing the lessons learned and highlighting opportunities for further exploration, we aim to foster continued innovation and advancement in this rapidly evolving field.

## 4.1   Summary

Throughout this paper, we have delved into the realm of Human-Machine Image Based Dialogs, exploring the fusion of computer vision and natural language processing to enable interactive dialog systems grounded in visual context. We began by defining the concept of Visual Dialog, elucidating their significance in bridging the gap between visual perception and language understanding. Motivated by the need to address accessibility challenges, advance AI technology, and enhance user experience, we embarked on a journey to explore the potential of Visual Dialog systems.

Our objectives were clear: to define the Visual Dialog task, introduce evaluation metrics, identify challenges faced by researchers, and highlight real-life applications. To achieve these goals, we meticulously examined the methodology behind Visual Dialog systems, encompassing dataset preparation, model architecture, training procedures, and testing methodologies. Through empirical evaluation, we revealed the performance of Visual Dialog models on VisDial and VQA datasets, showcasing their ability to generate contextually rich responses in dialog contexts.

However, our exploration also unearthed several challenges inherent in the development and deployment of Visual Dialog systems. From managing model complexity to addressing biases and ensuring fairness, and navigating real-world deployment challenges, each obstacle presents an opportunity for innovation and refinement.

In summary, our journey into Human-Machine Image Based Dialog has provided valuable insights into the capabilities, challenges, and potential applications of Visual Dialog systems. By understanding the nuances of this burgeoning field, we lay the groundwork for future research endeavors aimed at creating more robust, inclusive, and impactful dialog systems for human-machine interaction.

## 4.2 Applications

The potential applications of Visual Dialog systems span across diverse domains, offering transformative solutions to various real-world challenges.

- Virtual Assistants Visual Dialog systems can empower virtual assistants by enhancing their ability to understand and respond to user queries in a more contextually relevant manner. By integrating visual understanding capabilities, virtual assistants can assist users with tasks such as identifying objects in images, providing visual explanations, and guiding users through visual content, thereby enriching the user experience and increasing the efficiency of virtual assistant interactions.

- Educational Tools In the realm of education, Visual Dialog systems hold immense promise as interactive learning aids. By enabling dynamic dialogues between students and virtual tutors, these systems can facilitate personalized learning experiences, clarify complex concepts through visual demonstrations, and engage students in interactive problem-solving activities. Visual Dialog systems can also support remote learning initiatives by providing visual explanations and interactive feedback to students in online educational platforms.

- E-Commerce Visual Dialog systems have the potential to revolutionize the e-commerce experience by enabling more intuitive and personalized shopping experiences. By integrating visual search capabilities into e-commerce platforms, users can interact with product catalogs using natural language queries and visual inputs, receive personalized product recommendations based on their preferences and browsing history, and engage in conversational interactions to gather additional product information or make purchase decisions.

- Healthcare In healthcare, Visual Dialog systems can serve as valuable tools for medical imaging analysis, patient education, and remote healthcare delivery. By leveraging computer vision techniques, these systems can assist healthcare professionals in analyzing medical images, detecting abnormalities, and providing diagnostic insights. Additionally, Visual Dialog systems can facilitate patient education

by generating interactive visual explanations of medical procedures, treatment options, and healthcare information, empowering patients to make informed decisions about their health.

# Bibliography

[1] A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. Moura, D. Parikh, and D. Batra, "Visual dialog," 2017. visualdialog.org.