

Visual Question Answering using Transformer Architectures: Applying Transformer Models to Improve Performance in VQA Tasks

Vedant Singh

USA

ABSTRACT

Visual Question Answering (VQA) programs are an area of subject and field in computer science that seeks to develop technologies that enable the client to answer questions based on displayed images. Bert, Vision Transformers (ViTs), and multimodal transformers have aided the VQA systems significantly in understanding the relations between vision and text data. In this paper, these architectures are considered in relation to the scalability, dynamic attention mechanism, and multimodal pre-trained models of the prior CN-RNN hybrid models' weaknesses. Feedback activities connected with assistive technologies, healthcare, retail, self-driving vehicles, and creative industries indicate how VQA might be easily introduced and provide examples of likely relative positive and negative societal impacts such as bias, privacy, and inclusion. VQA systems are slowly becoming paramount for improving accessibility solutions, for instance, where a visually impaired person talks to the system to explain what a picture is about. However, VQA systems based on current transformers have some issues, notably from the point of view of computational complexity and reasoning capability. This paper covers the current state of research, issues, and the direction of further development, going further and noting that more attention should be paid to lightweight models, datasets from multiple domains, as well as the integration of human-generated data with AI. Accordingly, the identified results show that VQA systems can become one of the elements of context-aware, inquiry-based solutions for advanced applications in various fields.

*Corresponding author

Vedant Singh, USA.

Received: March 01, 2022; **Accepted:** March 08, 2022; **Published:** March 31, 2022

Keywords: Visual Question Answering, Transformer Models, Vision Transformers, Multimodal AI, Natural Language Processing, Deep Learning, Assistive Technologies, Ethical AI, Human-AI Collaboration, Lightweight Transformers

Introduction

Visual Question Answering, or VQA, is a highly interesting intersection between computer vision and natural language processing, which, in fact, involves creating machines capable of interpreting questions related to images to provide accurate answers to the questions. Described by some as a breakthrough in artificial intelligence, VQA poses challenging requirements for understanding real-world images and the contexts to ask questions about them. In the past, AI models shared features using Convolutional Neural Networks (CNNs) for Figure 1 visual outputs and Recurrent Neural Networks (RNNs) for textual inputs. Such approaches, however, were initially unable to work well in modality integration and did not perform so well in VQA's complicated tasks.

Deep learning of transformer architectures introduced an era of unprecedented AI advantage. Transformers brought dynamic attention modules that can pay attention to the specific parts of the input while preserving temporal relations. First coined for NLP, transformers rose to prominence for multimodal tasks because of the property of parallel computation and expansibility. This has been a significant improvement for VQA systems, as the models learned to fuse the visual and textual modalities more holistically. This paper discusses how to apply transformer-based architectures in VQA and the rationale for doing so based on the shortcomings of prior platforms. The suitability of transformers

for multimodal learning is considered, with further discussion on dynamic attention and self-supervision. Additionally, they are compared on datasets such as VQA v2 and GQA by assessing how effectively these models address distortions in real-world data.

Aside from performance, transformer architectures create new opportunities concerning user-oriented aspects of VQA. These system possibilities have complex implications in subject areas, such as the accessibility of assistive technologies to the visually impaired and user-interactive education. At the same time, using smart services and products entails a set of ethical challenges that must be solved to provide an essential level of fairness and privacy, namely data protection and the prevention of bias. Subsequent sections of this paper detail the development of VQA frameworks, explain transformer models and their modality extensions and assess the use of transformers for different practical issues. This paper aims to systematically review the existing approaches and advance a more comprehensive understanding of the opportunities and challenges associated with transformer-based VQA systems for future development.

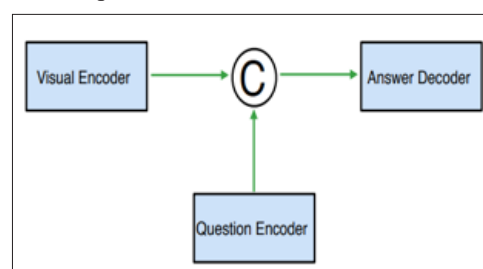


Figure 1: Flowchart of a Visual Question Answering Model

Background and Related Work

Overview of VQA

Visual Question Answering (VQA) is a cross-over task between computer vision, artificial intelligence, and natural language processing [1]. To answer questions correctly, an AI system must comfortably understand general information about an image and details related to a picture. Compared to other vision tasks like object detection or segmentation, VQA requires a higher level of understanding of vision components and their connection and relation to a given question. This makes VQA a difficult place to study the interaction of vision and language understanding across AI systems. For example, however, responding to a question such as "What color is the car beside the tree?" involves identifying objects (car, tree), pointing out the relations of objects, and linking all these to the query.

One of the most significant sources in further redefining VQA has been the emergence of datasets such as VQA v2 and CLEVR. These benchmarks are intentionally constructed to calibrate logical inference, contextual relationship, and photo-text correspondence viewing in particular models. While VQA v2 corresponds to real-life image backgrounds to questions, CLEVR presents well-defined synthetic environments to question and reason about. These datasets compel the researchers to think beyond VQA models to

apply the right models with a primary focus on reasoning instead of low-level pattern matching. One big issue in VQA is ambiguity and how the query depends on the context. Factual questions, for instance, "What is this?" can give simple answers, whereas it is not the same with questions that ask for relations between more than one object or for abstract knowledge, which makes the exercise much more challenging. This shows that besides providing relevant information to users, models also require rich contextual awareness and reasoning to address different types of queries the users pose.

One can apply principles behind VQA in various fields that have nothing to do with academic performance [2]. One of the uses of interpreting visuals is that it can answer questions. Since making contextual query interpretations vital, such as assistive technologies, interpreting visuals has potential uses in the following: This makes the task not just a technical one but also a move to develop AI systems that can better interface with reality. The progress of VQA research aligns with the general trend of development in AI, including the transition to end-to-end models and beyond. This is especially the case with the rising popularity of transformer-based methods, which also presents new opportunities for overcoming the shortcomings of CNN-RNN hybrids and expanding the range of potential applications of VQA systems.

Table 1: Comparison of VQA Datasets

Dataset	Focus Area	Key Features	Challenges Addressed	Example Question
VQA v2	Real-world images	Natural contexts, open-ended questions	Ambiguity, reasoning in real-world scenarios	"What is the person doing in the image?"
CLEVR	Synthetic environments	Controlled complexities, logical inferences	Spatial reasoning, multi-step reasoning	"Are there more red cubes than blue?"
COCO	Image-captioning	Object detection and caption generation	Object recognition	"What objects are present in the image?"

Transformer Architectures

Transformers have replaced attention-based mechanisms in the AI domain, allowing it to model dependencies within data modalities [3]. Initially created for NLP, models such as BERT have proven the applicability of attention in handling dependencies between tokens, boosting breakthroughs in text mining techniques, including class realization, summarization, or translation. Since they can handle inputs in parallel and are scalable, the computer vision tasks adopted them, leading to Vision Transformers (ViTs) and multimodal models for VQA. Like Vision Transformers (ViTs), researchers modified the transformer framework to work with image data by patching images into substrings or tokens. These tokens are then passed through the Self-Attention layer to capture the relations between different tokens so that local and even global visions can be captured. While getting hierarchical features of the image CNNs might be limited, ViTs efficiently maintain self-attention to attend to the relative parts of the image and thus are more flexible for capturing contextual relationships. Other multimodal transformers, including LXMERT and UNITER, expand these principles to encompass visual and textual input. They employ cross-attention mechanisms to synergize features from the visual inputs, bearing in mind that ViTs and text embeddings are from language transformers such as BERT. This approach helped the multimodal transformers perform the complex VQA tasks by notifying semantically richer image regions about the query.

be adapted to train for specific tasks as in the VQA v2. Pretraining increases performance while simultaneously ensuring that models can generalize across different domains, an unprecedented feature useful in real-world applications. Transformers have some limitations, mostly connected with their efficiency in terms of computation and scalability [4]. Standard transformers impose heavy requirements on their clients in terms of the necessary hardware, so they are not easily available. Several efficient transformer variants, including Linformer and Performer, have been proposed to deal with this. These developments make transformers progressively applicable for VQA and other multimodal tasks.

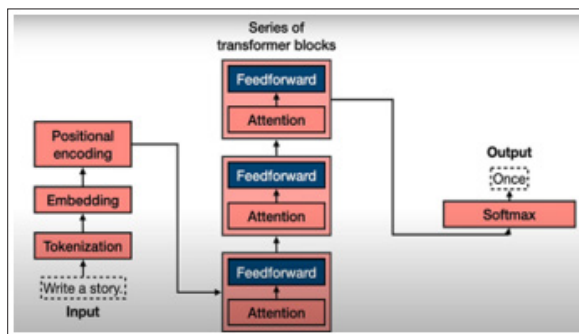


Figure 2: Simplified Architecture of Transformer

It has been very effective in utilizing the transformer architectures for the VQA. These models have been trained on data made of image-text pairs from real-world scenarios such as Conceptual Captions and COCO, which are general representations that can

Previous Approaches in VQA
 Before the transformer for VQA, CNN-RNN combined most of the models to process the visuals and texts in the question [5]. CNNs were used to determine the visual characteristics of

images, while RNNs or some of its flavors, such as LSTMs, worked on the textual representations. These models usually had attention mechanisms to pick out features or areas of the image related to the question. However, these architectures were brittle or generalization-deficient, especially when they tried to handle multiple datasets or perform smart searches.

VQA activities and the reinforcement process were presented as a method of improving the quality of the prediction. In another study, researchers trained their models to look for features in images and questions relevant to answering them based on the reward signals obtained from the passage. However, previous work based on reinforcement learning was computationally time-consuming, and fine-tuning was needed most of the time to get the best results. Some of these problems were tackled in one way or another in the latter hybrid models that included hand-crafted features for instance, giving models semantic segmentation or object detection outputs provided more information. However, these methods added a lot of complexity to the pipeline, and in most cases, the approach failed to address the issue of knowledge fusion between text and vision.

The application of attention mechanisms was a clear breakout in VQA research. As a result, attention-based CNN-RNN models were able to show a boost in accuracy by allowing models to prioritize regions of interest in the image based on a specific question. However, these models were still limited by the sequential data processing and were thus not as efficient and scalable as transformer methods. Since then, transformer-based architectures have evolved into the state-of-the-art for the VQA problem, providing end-to-end solutions for visual and textual information processing [6]. Their flexible long-range dependencies and dynamic modulation of different modalities eliminate many of the issues in the previous approaches and provide a new state of the ART with high performance and scalability.

Architectural Innovations in VQA using Transformers Unified Representations with Multimodal Transformers

Multimodal transformers have revolutionized how vision and text modalities are combined for VQA tasks due to the incorporation of cross-attention layers. These layers enable the models to dynamically map visual features of images with text embeddings developed from the language encoders. For example, Vision Transformers (ViTs) divide images into patch tokens, where each token is a portion of the image input, and textual data is tokenized and processed using language transformers, such as BERT. The cross-attention operation guarantees that visual and textual data are matched up in the correct context so that the model can learn about correlations between the data modalities.

This alignment process is very important to get the right answers. Multimodal transformers centered on semantically fluent areas of an image make fewer mistakes prevalent in prior models that treated all pixels with the same importance irrespective of their connection to the query. For instance, when the car next to the tree is posed, what color is it? The model can focus its attention on sections such as the car and the tree while possibly completely discarding sections that are not useful, such as people present in the image. Such sheer attention increases the precision of responses and the ability to explain how the model arrives at a given conclusion [7].

This is one of the strengths of multimodal transformers, as they can process inputs in different formats. These models can effectively

select and align features in the input regardless of whether the input is as detailed as a photograph or as simple as a black-and-white drawing to answer the question. This versatility of operation is useful for various applications, including instruction aids and highly technical operations. The success of multimodal transformers also reflects their scalability argument. Unlike previous models, which used two parallel pipelines for imagery and text data, transformers are a generalizable architecture that can be adapted to work on any dataset and multiple downstream tasks. This helps reduce the development cycle and makes extending the created VQA systems to other domains easier. While they are considered rather effective, multimodal transformers are not without their troubles [8]. The unlimited demand for computation and rich training data annotated by humans are still the major challenges. Still, recent developments in hardware and pretraining prevent these difficulties and, therefore, prompt multimodal transformers for VQA development.

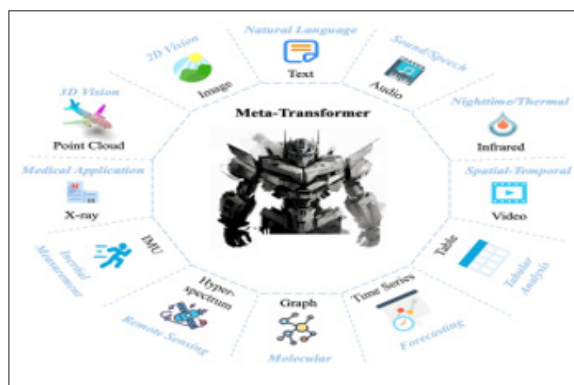


Figure 3: Meta-Transformer

Vision-Language Pretraining

Transformer-based VQA approaches that leverage and fine-tune the rich source and target corpora are now standard best practices. This process trains models with datasets like Conceptual Captions, which associates images with captions, or COCO, a very common dataset for image-captioning tasks. Most of these models thus take a pretraining approach where models are trained to learn mappings between the visual features and corresponding textual descriptions to lay down basic knowledge that can be fine-tuned for specific modes of VQA [9]. Another advantage of vision-language pretraining is that it enhances generalization. As in other diverse domains, such as VQA v2 and GQA, the robust benchmarks indicate that models VisualBERT and OSCAR achieve higher precision when fine-tuning with little specific task-related data. This is due to the fact that pretraining of these models involves optimization with respect to matching visual objects and textual description, thereby endowing them with rich multimodal relationships while answering diverse domains and question types.

Pretraining also helps to solve the problem of data deficiency, which is traditional in VQA investigations [10]. Data acquisition and selection of specific annotation tasks could be time- and cost-consuming, especially when additional domain knowledge is needed, example, in the medical field. Vision-language pretraining solves this problem by using public datasets to learn general features, of which the model is then fine-tuned on specific datasets for special task requirements. Pretraining serves to attend to the contingencies surrounding individuals more than advancing performance, and in doing so, promotes more coherent and understandable VQA systems.

Through intermediate representations of images and words and achieving correspondence between the objects in the images and their textual descriptions, the pre-trained models explain their workings. For example, when answering a question about the image, the model can point out which regions and textual features it found useful while answering the question, which creates very important transparency, especially in fields such as medicine and security. There are certain drawbacks of vision-language pretraining. This survey shows that data quality matters a lot in the performance and bias of the model during the preprocessing stages of the pretraining dataset. Broader data sets that are not inclusive can result in ineffective models for underrepresented situations. Preprocessing the datasets and formulating a better pretraining approach solve these prejudices.

Hierarchical Reasoning Mechanisms

A major test case where hierarchy plays an important role is one of the most significant problems connected to VQA – those that cannot be solved by simple recognition of objects. CLEVR-style datasets have been created to measure different logical inferences, spatial relations, and attribute comparisons in the scene. For example, respond to questions like "Are there more red spheres than blue cubes"? Expect the model to recognize objects or segments of the scene, categorize some attributes, and count objects while keeping track of the rest of the scene [11].

Because these tasks require hierarchical attention mechanisms, Transformers shine in them. Such features have been achieved by stacking many self-attention layers, thereby enabling these models to progressively capture several levels of interactions between entities in an image. In the earlier layers, the model may only concentrate on getting simple features, such as shapes and colors. Relatively more complex queries are answered as the data moves through other layers, and the model learns to integrate these extracted features to provide solutions. The other advantage of transformers in hierarchical reasoning is that they can solve a sequence of reasoning problems. Traditional models fail in scenarios that involve handling information based on a step-by-step process, such as "What is the shape of the object to the left of the red sphere?" It is evident how transformers can model these dependencies with their global attention mechanism while providing accuracy and context awareness in their reasoning.

Another advantage stems from the hierarchical nature of transformers, from the domain layer to the application layer. These capabilities also harness the transformers' flexibility in dynamic scenarios. For real use cases that involve robotics or autonomous systems, where a prompt answer to a query based on contexts is necessary, transformers can easily learn from new data coming into the system. This aspect makes them differ from other models that are rigid and generic and developed to meet specific tasks only. As the results highlight, hierarchical reasoning still presents a difficulty concerning transformer-based VQA systems. Arguments requiring interpretation, evaluation, etc., such as "What do you think of this scene regarding the nature of the character's motives?" are still challenging for contemporary models. Mitigating the above will also demand the development of model architecture and other training procedures [12].

Table 2: Hierarchical Reasoning Capabilities

Reasoning Level	Task Example	Transformer Advantage
Basic Feature Extraction	Identifying shapes and colors	Focus on local features using self-attention
Intermediate Object Mapping	Object relationships in a scene	Cross-modality alignment for contextual reasoning
Advanced Logical Inference	Solving multi-step spatial problems	Long-range dependencies through hierarchical layers

Transformers for Efficient VQA

The high computational complexity of transformers has always been a key issue in AI, and it especially affected the study of high-consuming applications such as VQA. Standard transformers consume large amounts of memory and computational resources. Thus, it is challenging for researchers and practitioners facing restricted infrastructure requirements. To this end, Linformer, Performer, and other improved transformer models have been proposed to control the time complexity. Efficiency is maintained by mimicking the self-attention mechanism while scaling the computational cost for this operation from quadratic to linear in the input size of Linformer. This approach allows fine-tuning and executing transformer models on comparatively low-powered devices like smartphones or edge devices. Likewise, Performer innovates kernel-based approximations of self-attention and only enhances efficiency while retaining this architecture's reasoning ability.

More efficient transformers tremendously impact the size and expansibility of VQA systems. Such models decrease the necessary amount of hardware to achieve real-life applications of VQA solutions, e.g., for visually impaired people or robotic systems in the industry. This democratization of technology is helpful in making sure that the positive outcome of employing VQA can become available to as many spectrums as necessary [13]. The use of efficient transformers makes real-time applications possible. In burgeoning areas such as self-driving or live gaming, where time availability is paramount, these models can sort out queries and form answers quickly, paving the way for better UX and improving operational limits. Efficient transformers still have challenges where high performance is a bottleneck for a higher computational cost. Although they distribute the accumulated information to require less memory, their efficiency may be lower than regular transformers, depending on how deep the problem understanding is. In these trade-offs, more future work will have to be done to zero in on the best architectures for maximum efficiency and effectiveness in VQA.

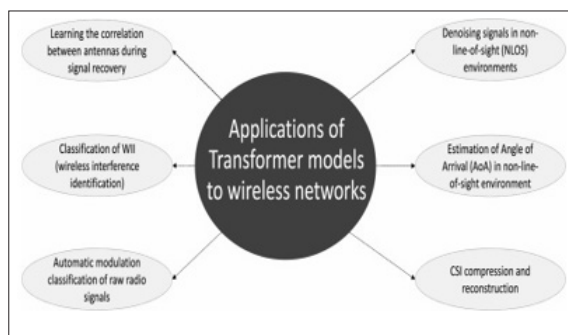


Figure 4: Generative Pre-Trained Transformer

Applications and Practical Implications of VQA using Transformers

Table 3: Applications of VQA in different Domains

Domain	Application	Example Query	Key Benefits
Assistive Technologies	Scene navigation, object recognition	"What is on the table?"	Independence for visually impaired users
Healthcare	Medical image analysis	"Are there fractures in this X-ray?"	Faster and accurate diagnostics
Retail and E-commerce	Visual search, automated FAQs	"What material is this sofa made of?"	Enhanced customer experience
Autonomous Systems	Navigation, object detection	"What obstacles are ahead?"	Improved safety and operational efficiency
Creative Industries	Interactive gaming, scene analysis	"Who is speaking in this movie scene?"	Inclusive and immersive user experiences

Assistive Technologies

Transformed-based VQA systems have revolutionized assistive technologies by allowing visually impaired persons to interact with their environment in real-time [14]. These systems use multimodal learning to contain scene descriptions, object recognition, and contextual questions. For instance, a person could speak to the system like, 'Tell me what's in front of the table.' and obtain a precise description of the items, their color, and even their position in space. This capability promotes enhanced independence, enabling navigating through contexts and making decisions independently.

Scene navigation is an example of the major use of VQA since the systems can offer directions as they point out relevant objects in the scenes. Questions such as "Where is the exit?" or "Can I see an obstacle in front of me?" can be answered accurately and assist a user in cases where a person seems to be lost in unknown territory. This has important implications for accessibility in public areas, as these problems are still present for persons with visual impairment. In addition, these systems, together with the Navigator, can support an individual reading and comprehending texts in the surrounding environment. For instance, a VQA system can answer questions about the text on signs, documents, or electronic display surfaces. This feature is rather valuable as it has no equivalent in tactile or audio format, like audiobooks or braille. It is even more useful when the VQA includes a real-time object detection component with context-based answering. For example, visually impaired users may ask: "What is the person in front of me doing?" and receive a description of actions and facial expressions. This is because it develops social interaction since people can learn different gestures in a conversation better [15].

Issues hinder the implementation of these systems so that they are cheaper and available to everyone. Constraints associated with the use of equipment and space, the size of the model, and the requirement for an internet connection can hamper its implementation, especially in low-resource environments. New research should involve the investigation of novel methods of achieving affordable and lighter solutions that would make these transformative technologies easily accessible to the relevant communities [16].

Healthcare Applications

They have the futuristic prospects to revolutionize the diagnostic process in healthcare to analyze medical images using transformer-based VQA systems. It can generate responses to questions about aberrations observed from X-ray, CT, or MRI scans to help radiologists determine from thicker scans. For instance, a radiologist could ask the system, "Are there any fractures in this X-ray?" and get a precise answer pointing to the area of the picture containing the answer to the given question. This capability decreases diagnostic mistakes and enhances the speed of the decision-making processes.

VQA is also important in educating and training medical personnel and health care professionals. Using VQA and other interactive learning tools may provide medical schools and hospitals the means to work through examining complex medical pictures with medical students and residents. To build such an understanding, examples of questions that can be asked include "What abnormalities are likely to be seen in this MRI scan?" when several students discuss their case, the tutor is also able to explain in detail together with illustrations that can lead to a broader understanding on structures of the body and diseases. These systems help educate patients; doctors can describe certain diseases or potential therapies using images [17]. Patients may wonder, "What does this part of the X-ray mean?" and get a brief explanation with illustrations and annotations of their sub-images. This increases interaction between medical practitioners and patients to ensure that most of the information being passed across is understood and believed by the patients.

In surgical planning and intraoperative assistance, VQA systems can be valuable in making decisions based on medical imagery. For example, in performing minimally invasive surgery, the surgeon could ask the system for some orientation, such as anatomical landmarks or the location of certain structures. This helps minimize cognitive loads and also adds precision, especially in areas that are complicated or risky, such as surgery. Incorporating VQA in healthcare delivery prompts crucial questions regarding data privacy, model accuracy, and the importance of versioning, which requires domain-specific training. Therefore, for these systems to be usable in clinical settings, these systems have to be trained on diverse datasets, and most importantly, all systems must meet the set legal provisions. Overcoming these challenges should go a long way toward helping realize VQA's full potential in the health sector.



Figure 5: The Diagnostic Revolution: How AI is Transforming Medical Imaging

Retail and E-Commerce

Applying VQA technology in e-commerce platforms is rapidly growing due to applications like visual search and automatic

FAQs. Customers can order specific product details, such as "What material has been used to make this sofa?" or "Is this dress allowed to be put in the washing machine?" Image understanding is performed by using transformers that extract the relevant characteristics of the product images and the metadata to give accurate contextual answers that assist customers in purchasing a particular product [18]. VQA systems also enhance the convenience of product finding, especially when using the visual search process. For instance, customers can scan an image of a preferred product and query, "Can you supply similar items?" The system can establish the image and push out the essential characteristics such as color and design, then suggest products to be purchased. The design concept of this unique shopping service is inherently friendly and loyal to customers. Besides the applications in customers' interfaces, VQA integrates posterior processes, including supply chain management. A sample of the questions that warehouse staff can use to query systems include "How many of these products are in this warehouse?" or "Do you find any flaw in this batch of products?" Such systems integrate visual inspection and automated answering, improving the working process and minimizing human interference.

Retailers can also utilize VQA for marketing and customer relations. For instance, the current product information on VQA interactive displays can enable users to drill into product specifications in real time. These kinds of questions: "What are the dimensions of this table?" or "products that can interconnect with this device?" product information can be provided in a moment, improving the buying process. When VQA is implemented in retail and e-commerce businesses, an array of issues need to be solved: First, to achieve high accuracy of the answers across a wide range of product categories. Second, to handle the computations required to provide immediate responses. Future trends ought to incorporate a specific e-commerce VQA model type for retail and enhance compatibility with present e-commerce applications.

Autonomous Systems

Transformer-based VQA systems are revolutionary for autonomous systems, allowing robots and other automatic tools to understand the environment and answer questions [19]. These systems enhance mobility since robots can now answer questions such as: "Where is the blue box?" or "What type of obstacle confronts me?" Thus, this capability improves the working ability of robots in some industries, such as logistics, and manufacturing. VQassist Automated warehouses in tasks such as identifying misplaced items, detecting lost inventory, and improving the arrangement. For instance, a VQA model robot can walk to an aisle and, in unison, respond, "Are there any products missing from this section?" These systems include a feedback recipient that sorts data and elements of reasoning that lower the mistakes and enhance business operations in comprehensively intricate tasks.

VQA is also very important in driving self-driving cars since the environment must be interpreted properly for safety's sake. For example, an automobile could use a VQA system to understand road signs, detect an accident-prone area, or ask a question such as, "Is the way clear?" This improves decision-making processes and makes travel safer when operating in uncertain conditions [20]. In robotics, VQA allows users to ask questions about the tasks the robots are performing or the environment. Questions such as "Where is the package you are delivering at the moment?" or "What's on the table?" may be answered satisfactorily, thus increasing interactivity between humans and robots in industrial and service environments. Integrating VQA systems in autonomous systems presents challenges concerning the system's real-time response and scalability. Supporting these models' capability to run on edge devices while maintaining the

same level of accuracy is a key avenue for further work. Overcoming these challenges will inevitably broaden the opportunities for using VQA in autonomous systems.

Recreation and Performing Arts

The entertainment and creative industries use VQA technology and tools to make their products more relevant and accessible. In interactive gaming, VQA systems enable players to ask questions about visible parts of the game environment, making it more realistic. For instance, players may ask questions such as: "What is the goal of this mission?" or "What is it I can use?" These dynamic interactions make games more enjoyable and active, improving the general experience. The VQA systems' tasks in media production include scene analysis and content description. For viewers with vision impairments, these systems describe scenes in movies or on the TV screen, using questions like, "What is going on in the scene?" or "Who is speaking?" It also enhances the culture of diversity to make media content properly reachable to all interested people.

Creative workers apply VQA technology in the movie industry to enhance their productivity in innovative areas such as animation and VFX. For example, the designers can ask systems about some aspects of a scene, such as "Which assets require resizing?" or "How is it lit in this frame?" The routine questions are eliminated as they are run through the VQA systems, freeing the professional's time for creative input alone. In marketing communication, VQA systems allow brands to build advertising campaigns that give consumers a more participative experience. For instance, an advert for a product could have a section where a viewer can ask questions such as: "What is in this drink?" To ask the questions 'How' or "What makes this product eco-friendly?" Such a level of interaction helps promote stronger engagement with the target audiences [21]. Implementing VQA in the entertainment and creative industries is not easy, and there are several drawbacks, such as the need to perform in real time and the flexibility in the type of content. For these systems to apply in these fields, it will be crucial for them to be able to meet the specific requirements of creative processes and be accurate at the same time.

Societal Impacts of VQA Technologies

There is much at stake when applying Visual Question Answering, sometimes called VQA, to everyday use. Far from simply improving accessibility and knowledge, these systems pose ethical questions that cannot go unasked. As suggested earlier, VQA has the capability of helping oppressed groups and promoting equity, in addition to changing business sectors. However, problems concerning data protection, prejudiciality, and the proper application of artificial intelligence are still significant. To this end, this section looks into the effects of VQA on society in terms of accessibility, education, ethics, and future [22].

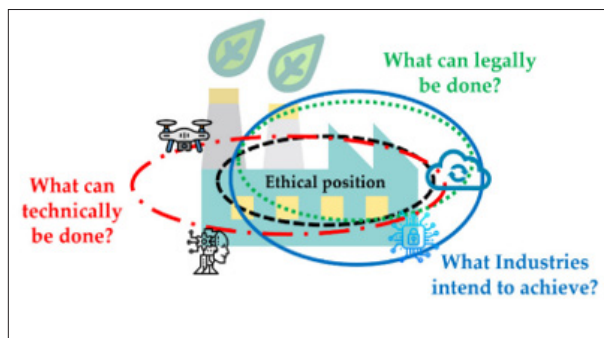


Figure 6: Ethical Dilemmas in Enabling Technologies used in Industrial IoT (Diagram Adapted from IBM Model for Ethical Analysis)

Accessibility and Inclusiveness

Another important social contribution of VQA technologies is that they may be used to improve the availability of services to people with disabilities [23]. Focusing more on visually impaired users, they enhance their real-time interaction with the immediate environment. For instance, they can explain environments, objects, and texts, thereby enabling users to read space and text and move on their own. It also leads to a higher level of independence and has a positive impact on the quality of life of people with low vision. VQA can also benefit people with cognitive issues by converting intricate visual data into plain answers they can comprehend. For instance, the following questions are real: "What does this chart depict?" or "What do you see written on this sign?" and get brief and to-the-point appropriate answers. This capability makes information easily available to the public and can help eliminate certain parameters that have kept certain groups of people away from certain information.

Incorporating VQA into public projects can also enhance the population's participation agenda. This means that where there are VQA-enabled kiosks or mobile applications, users with disability should be able to get detailed directions on how they can get the services or access this facility on their own. For example, questions such as "Where is the nearest elevator?" or "How to use this machine?" can be answered quickly and have precise responses to the questions. The overall user interface experience will be pleasant. Although these are true, some problems regarding implementing technologies in VQA and the distribution of their corresponding advantages still exist. They include high cost, accessibility in underdeveloped regions, technological advancement hitherto observed in developed nations, and the fact that there may be a language barrier in how they are programmed or designed. An attempt at creating cheap and cross-lingual VQA systems is crucial to extending the approach's benefits to as many people as possible [24].

Changing Education and Learning

Applying VQA technologies can change educational activities and make them more communicative and interesting. In classrooms, students can use VQA systems to interact with image content more interactively by posing questions such as "What does this part of the diagram mean?" or "What did this map tell?" This approach fosters curiosity and enhances learning, especially for topics involving topics of graphic data, such as geography, biology, and history. Speaking of specific benefits for educators, VQA offers tools that will help develop an individual approach. By asking questions about the visual content, teachers can encourage their students to use their potential in problem-solving and critical analysis of stimuli. Professionally, in virtual learning situations, VQA systems can perform the roles of virtual tutors to respond to students' questions on any material shown, such as slides or even videos, in real time.

VQA also plays a very important role in overcoming barriers to education [25]. These systems can read the content to blind or Dyslexic students or translate text into graphics or explanations, thus engaging the challenged students in the visuals as well. This is true for students in rural or other hard-to-reach regions, where VQA-enabled devices can enrich the teaching process. The use of VQA in education has some problems, such as the problem that results from the divide in the use of technology resources. Schools in poor areas might be unable to afford the technology, which would further widen the gap. These tools must also be available for students of all SES statuses since policy and education should not only support those in the upper SES.

Ethical Issues Regarding the use of VQA

As with many emerging technologies, implementing VQA technologies raises profound ethical concerns in data privacy, bias, and fairness. These systems use big data as a training dataset, most of which can be considered unique or contain personal information. Data management is vital to user privacy and the appropriate use of VQA applications; hence, it should be implemented and followed to the letter. Another important issue in VQA systems is bias. In this work, bias is defined as the ability of the model to answer questions in a way that may not be fair to all the participants in the scene [26]. Those trained on some datasets will have certain biases or provide unfair or inaccurate responses. For example, VQA systems trained mostly in featured areas or ethnicities may fail to understand or correlate visuals to underrepresented cultures. To overcome this problem, actions should be taken to diversify the training datasets and dourness audits in the model development.

Transparency can be said to be critical in minimizing ethical risks in the deployment of VQA. People should be able to realize how VQA systems arrived at their answers, especially when used in sensitive areas like medical or law enforcement. Making the activities of the models more transparent can promote accountability and trust in their decisions. Ethics concerns the use of VQA in surveillance or any other invasive capacity. Because these systems can analyze and interpret visual data quickly, the potential for their misuse tends to affect people's privacy and consent rights. It is, therefore, important to establish good and well-checked regulatory/ethical boundaries to counter these corruptions.



Figure 7: The Ethics of AI

Future Implications and Considerations

The further development of VQA technologies has to improve them and increase their applications in society. In both cases, the resources available for visual question-answering models must be kept at the bare minimum. This is where lightweight, low-cost models become important to avoid the obvious impossibility of implementing such systems in resource-constrained environments such as rural schools or low-income areas. These developments will allow the state to erode the digital divide and extend more VQA benefits to a larger population. Another objective is to include multilingual and culturally sensitive VQA systems as one of the urgent tasks. Modern models do not work well beyond the first languages and cultures that are part of most current AI systems of solutions. Thus, the display of stereotyped responses for localized images and introduced different languages, cultures, and datasets to make VQA systems more universal. One also envisages the possibility of expanding the sphere of application of VQA together with other AI technologies, for instance, generative models and robotics [27]. For instance, integrating VQA with conversational AI will allow for the development of systems that, besides answering questions about images, can actually have a conversation with users regarding visual content. Likewise, VQA-enabled robots might augment industry automation animal human tasks such as home care for older people.

While VQA technologies are progressively integrating, developing inter professional relations between technologists, policymakers, and ethicists will be vital. Those who decide on courses of action at the highest level of the organizations must set rules on how innovation can be encouraged while maintaining principles of social justice and privacy. Simultaneously, it is possible to have active discussions with ethicists and other community members to coordinate possible risks and determine whether specified VQA applications benefit the public. The actual benefits of VQA technologies to societies will also depend on their compactable and scalable advantages to society's existing ethics. By combining efforts to meet current difficulties and focusing on inclusive designs, VQA constitutes the prospect of a positive global impact on people and communities.

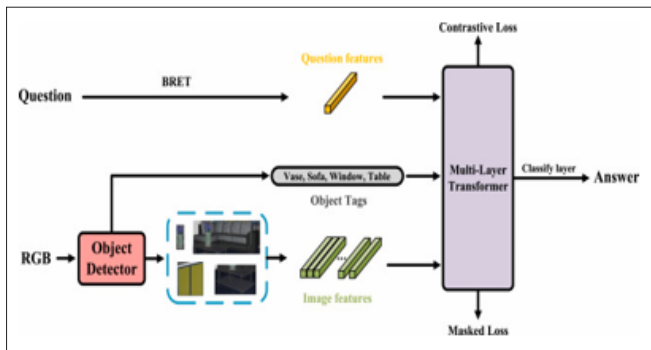


Figure 8: Transformer-Based Vision-Language Alignment for Robot Navigation and Question Answering

I/O Complexity and Hardware Requirements

Transformer-based systems for Visual Question Answering (VQA) are a giant step forward in AI systems, but they also come with issues [28]. These limitations are fundamentally due to computational complexity and size, as well as issues related to applying reasoning to subtle or vague queries. Another drawback of using transformer-based models used in VQA is their high computational complexity. Standard transformers, including those for shared domains, are memory and computationally intensive. This is because self-attention operations increase quadratically and, thus, are not scalable with input size. In VQA tasks, this problem is further aggravated by the need to process both image patches and tokens, further complicating the architecture and increasing the model size, thus making these models computationally expensive and hard to deploy on common hardware.

The dependence on high-performance GPUs or TPUs is another constraint because it can be impossible for researchers or organizations with small budgets to access them. Training transformer-style models take a large number of computational resources, which are expensive—something that can lock out researchers on a small budget or institutions with limited computational resources. This results in many people being unable to develop and implement high-quality VQA systems because of the lack of funding and knowledge needed to complete such a project. Several optimizations like Linformer and Performer introduced transformer versions with less dependency or demand for memory and time usage to address these limitations. As these approaches demonstrate reasonable applicability, they more frequently imply giving up either interpretability and training performance or logic reasoning to some extent. Efficiency versus performance is another focal point for future development since VQA systems are rapidly gaining importance as key resources in limited scenarios like education and healthcare.

Reasoning Disability and Vagueness in Questions

Transformer-based VQA models also exhibit limitations in reasoning tasks that demand high-level thinking or logical processes, chaining, and prior. Knowledge. For instance, probes such as “Can I predict what could happen next?” or “What kind of feelings does that person convey in the image?” require a degree of context processing and logic that current models struggle to accommodate. These tasks necessitate combining visual hints with knowledge about the world, which is challenging for models built mainly for data sets without motion [29]. Another major obstacle is that queries themselves can be quite ambiguous. It indicates that VQA systems may misinterpret many ‘fuzzy’ or ‘contextual’ questions because they are based on learned patterns derived from data. For example, a question such as “What do you think that is?” can search for different results, which may differ depending on the question and the picture itself. However, transformers, as powerful as they are, cannot actively ask for clarification, let alone learn how to handle dynamic changes in the inputs provided to them that may lead them to output errors or completely irrelevant results.

Some current methods to solve these reasoning difficulties are pretraining on various massive corpora and using external knowledge. However, these approaches are not free from problems as they fundamentally rely on the quality of training data that is available. Preconceptions can be especially magnified when it comes to data sets, which might lead to reasoning faults or the impossibility of applying the data to different cases. Awakening from these reasoning limitations will need a new model architecture that is more sensitive to context and capable of learning from it.

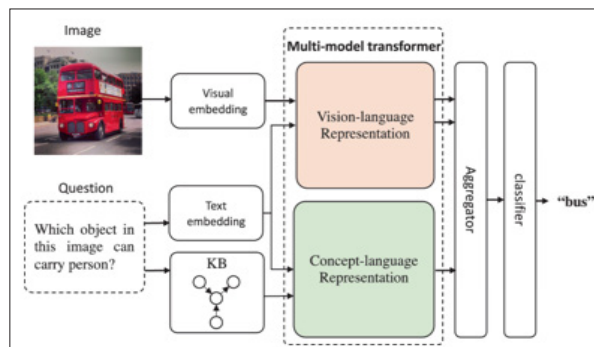


Figure 9: Knowledge-Based VQA

Future Directions in VQA Research

The development of Visual Question Answering (VQA) systems is one of the exciting Research topics in AI today, especially focusing on addressing existing shortcomings and extending the applicability of VQA systems. To enhance the performance of VQA, the researchers need to consider the following aspects: low-weight

- Low-weight model
- Improved reasoning mechanism
- Variety of datasets
- Mutual interaction between human and deep learning models

These directions will define the future development of VQA systems and make them more open, usable in many domains, and socially beneficial.

Table 4: Future Directions for VQA

Research Area	Focus	Expected Outcome
Lightweight Models	Optimization for resource-constrained devices	Increased accessibility and efficiency
Enhanced Reasoning	Better handling of abstract and multi-step tasks	Improved accuracy and versatility
Diverse Datasets	Inclusion of culturally and geographically diverse images	Fair and inclusive AI systems
Human-AI Collaboration	Interactive and adaptive systems	Improved user engage

Lightweight Model Design

Deeper transformer architectures depend solely on resources; streamlining such models is significant for running VQA systems on edge devices such as smartphones, drones, and robots [30]. Modern VQA models are computationally expensive and have practical implications that hinder real life, especially in areas of low resources. This can be concurrent with introducing the lightweight transformer, which may help reduce memory and processing and thus enhance VQA system feasibility and usability. Pruning, quantization, and knowledge distillation are used to make transformer models more efficient to make transformer models more efficient. These methods help decrease the model dimension and can cause limited deterioration of the model's performance. For instance, knowledge distillation is a process that preserves the accuracy of the given data by moving the knowledge from large-size models to small ones while cutting on the computational burdens. These optimizations are crucial in deploying portable VQAportably, really embedded devices.

Another exciting line of research includes raising cross-over between transformers and more efficient approaches to machine learning. For instance, while convolutional neural networks often achieve efficient feature extraction results at the start of the process, the development of lightweight transformer modules to handle the reasoning aspect can be quite effective. These hybrid models are suitable where in-real processing is imperative, namely in robotics and augmented reality. The scalability of lightweight VQA models on edge devices also wins in terms of privacy and security feats. These models decrease the likelihood of a data breach and latency since the data is processed on the system, not the cloud. This is even more paramount, especially when working on application issues within areas of sensitivity, for example, health or user devices. For these advancements to be realized, efforts from AI scientists, hardware developers, and industrialists need to be combined. Optimizing models specifically for certain gadgets or operating systems, including smartphones or wearable ones, will allow the designing of lean VQA systems that satisfy as many customers and tasks as possible.

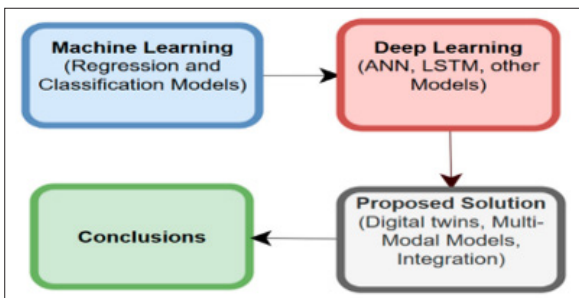


Figure 10: The Structure of the Paper Includes Machine Learning Models and Deep Learning Algorithms, and its Proposed Solutions.

Enhanced Working Memory

Improving the ability of VQA systems to reason is the critical area of development for facilitating agendas that involve more substantial and diverse tasks and practical use. It is not well equipped to perform pattern matching, abstract reasoning, putting things into context, and solving multi-step inferences that current models only attempt poorly. Enhancing the transformers to address these issues will greatly extend the spectrums of VQA applications. One is complexity modeling, in which models are required to realize contextual or figure-of-speech interpretations. For instance, responding to a question like 'What do you think it means about some people in this scene?' requires knowledge of feelings, plans, and culture. Current VQA systems have limitations in transferring image content to abstract reasoning. Hence, incorporating external knowledge bases and world models can enhance the VQA systems.

Another vital region for development is contextual reasoning. The VQA systems should be able to filter out the necessary and the unnecessary, alongside filtering and selecting the right context for the image query. Research suggests that the development of basic forms of attention leads to advancements such as dynamic attention, attention layers, and hierarchical attention, enhancing the model's focus on relevant features while disregarding irrelevant ones. Another urgent issue involves reasoning, where the system has to find the answer to multiple queries connected in terms of logic. Activities such as solving geometric problems or identifying a sequence of occurrences require sequential thinking strategies. Domain experts have identified new approaches, such as graphical representations and transformers with memory, to encode different relations in a superior way.

Advanced argumentation also challenges interpretability and the public's trust. Models that can describe how and why the particular outcomes fit the objective, such as pointing out the specific regions of an image or solving a problem in a sequence, will encourage more accountability. For instance, an app or location is very sensitive, such as in health or policing. Some of the reasoning difficulties listed above can be solved through the development of architectures to solve the problem and training techniques to improve the design. Thus, including diverse reasoning tasks during the training process and using the datasets of pretraining multimodal models will help ensure that current VQA systems are ready to address diverse real-world questions.

Multicultural and Multicultural Data

It is crucial to build diverse datasets to try to solve the VQA problem more fairly and for a wider range of situations [31]. In current datasets, issues arising from geographical, cultural, or demographic prejudices lead to inadequate models for handling underrepresented cases. Addressing these limitations will require enlarging the variation of training datasets and improving the VQA system's performance for users across the globe. The data set should be as diverse as the real-world scenarios that can be encountered, such as images from different cultures, locations, or situations. For example, training data should have pictures of rural and urban contexts compared with images containing diverse people and actions. This will assist VQA systems in extrapolating across different domains, including helping farmers in remote areas or helping them navigate within complex city environments.

Another important feature to consider is the presence of user-oriented contextual metadata. In most real-world scenarios, apart from the image contents, it is often necessary to extract

additional textual or contextual information from an image. For instance, radiology datasets need to include the patient's medical history as well as their clinical notes, which makes the answers given by the VQA system even more relevant. The race for increasingly diverse datasets also requires attention to data collection ethics. Guaranteeing that data samples used by AI algorithms are representative without violating privacy or consent remains highly important in creating socially beneficial AI [32]. Therefore, cooperative research by governments, universities, and local non-profit organizations can create ethical data collection and usage rules. When evaluating the advancement in this specific aspect, scholars have to employ indicators that capture equality or the lack thereof in VQA outcomes. Such benchmarks should assess how these models perform in various cultural, social, and linguistic situations and establish the working principle that inclusiveness should never be compromised in VQA.



Figure 11: Multi-Cultural Team Challenges and Solutions

Human-AI Collaboration

The future of VQA involves human and AI integration. VQA systems will work in synergy with humans in making decisions, and the outputs supplied by the VQA models will be easily understandable. Such systems should be seen as augmenting human knowledge and experience rather than substituting for it in some wholesale fashion in different areas of organizational functioning. One such avenue with growth potential is the incorporation of the ability to narrow down further or explain the original question using VQA systems. For example, if a system faces an ambiguous question, it could involve the user in a conversation to obtain more information about the question. This dynamic interaction makes the system better capable of giving precise and relevant answers to the users, enhancing satisfaction and confidence.

In the professional field, VQA systems can be useful as decision aids. For instance, in the health sector, the VQA system might be applied so that a doctor is likely to pose or prompt a question like What pathologies are present in this scan? Therefore, the system's response can be the doctor's second opinion, making decisions based on their findings. Likewise, in manufacturing, other engineers query robots with VQA to detect defects or optimize work; thus, humans and machines cooperate. It is worth adding that explainability is one important aspect of effective collaboration between humans and AI. The authors found that users require an indication of how VQA systems derive the answers, especially in critical scenarios. Items such as highlight bars, explanation paths, and confidence measures can help establish some accountability and let users judge the quality of the system's results.

Moving beyond vocational relevance, AI-enabled human Fig 3, for instance, can promote creativity and learning in VQA. For example, artists can paint using the VQA system as a tool to generate visual concepts. Students are also able to partake in meaningful learning

sessions. These systems can become collaborators in exploration, which helps generate ideas and cull information. To ensure that the incorporation of artificial intelligence can optimize human-AI interaction, the VQA systems developed by researchers must feature flexibility, aptitude for explanation, and human interfaces. These systems can facilitate users from various domains and applications through all the explored cases and analyzed concepts.

Conclusion

VQA acts as a platform incorporating vision and language understanding and transformer integration, one of the breakthroughs in artificial intelligence. These systems have been proven to have closed the gap between computer vision and natural language processing to provide complex interfaces for handling visual and textual information. Through self-attention, cross-modality, and hierarchical reasoning mechanisms, transformers have elevated the efficacy and applicability of VQA systems. Regarding the mentioned developments of the transformer-based VQA models, there is still a long way to go with multiple challenges and more improvement possibilities. Another significant shift attributed to VQA systems is the enhancement of real-world application possibilities in multiple domains. It spans applications that range from accessibility technologies for providing the differently-abled with a chance to be fully included in society to the healthcare sector, where VQA has the potential to help the right diagnosis be made within the shortest amount of time possible. It is also widely used in retail, e-commerce, autonomous systems, and entertainment, where the VQA systems play an important role in improving efficiency, interaction, and access. Such applications show how transformer-based VQA can be applied to various tasks, proving that such models open a new chapter in how people engage with visual/blob and textual/data content across their lives.

The use of VQA technologies in the Global environment is a significant issue that presents new key questions that must be solved for the technology to be valuable. Computational requirements and, hence, hardware costs still hinder these systems. Thus, they are largely available for companies and institutions that are financially capable. Also, there is a lack of effective reasoning within current VQA models – applying it to abstract queries and multi-step problems is difficult. Solving these challenges will require changes in the model design, training process, and design and construction of lightweight transformers for use in resource-limited environments. This is the power of VQA systems in society, but power has its costs and triggers accountability. Challenges of bias, data privacy, and benefit distribution impact appropriate data collection processes, training, good governance practices, and risk management. In addition, the relationships between people and AI in VQA can be strengthened to enhance the potential advantages as the AI system learns from human knowledge. In contrast, people benefit from speedy and accurate computations of the AI system. Maintaining a clear understanding of the logic and operation of such systems will play an important role in taking responsibility for their outcomes in fields that require special attention to privacy, such as medicine and police work. There are great opportunities for the further development of VQA research and its use in practice; these are closely connected with the responsibility for new technologies. If current drawbacks are resolved, the reasoning of VQA systems can be improved, and new models can be designed to be more accessible; VQA systems can help revolutionize different industries and aid people. It will be critical that the researchers work in close cooperation with policymakers and industry actors to ensure that the resulting systems are built and implemented in a way that is optimally beneficial to society and not detrimental.

With the future improvements in transformer-based VQA models, these models have the potential to be the fundamental component of smart and interactive systems in the increasingly miscellaneous digitalized humanity.

References

- Holland VM, Sams MR, Kaplan JD (2013) Intelligent language tutors: Theory shaping technology. Routledge.
- McCaslin M, Daniel TH (2013) Self-regulated learning and academic achievement: A Vygotskian view. In Self-regulated learning and academic achievement 213-238.
- Chen J, Xu H, Zhu L (2012) Internet of Things.
- Agarwal D, Long B, Traupman J, Xin D, Zhang L (2014) Laser: A scalable response prediction platform for online advertising. In Proceedings of the 7th ACM international conference on Web search and data mining 173-182.
- ANI, Narrow AI, General AI (2012) Optical Character Recognition (OCR).
- Yang BS, Zhou ZH, Gong Z, Zhang ML, Huang SJ (2014) Advances in knowledge discovery and data mining. In Proceedings 405.
- Healey C, Enns J (2011) Attention and visual memory in visualization and computer graphics. IEEE transactions on visualization and computer graphics 18: 1170-1188.
- Wagner J, André E, Jung F (2009) Smart sensor integration: A framework for multimodal emotion recognition in real-time. In 2009 3rd international conference on affective computing and intelligent interaction and workshops IEEE 1-8.
- Kumar A (2019) The convergence of predictive analytics in driving business intelligence and enhancing DevOps efficiency. International Journal of Computational Engineering and Management 6: 118-142.
- Deng L (2012) Three classes of deep learning architectures and their applications: a tutorial survey. APSIPA transactions on signal and information processing 57-58.
- Zia MZ, Stark M, Schiele B, Schindler K (2013) Detailed 3d representations for object recognition and modeling. IEEE transactions on pattern analysis and machine intelligence 35: 2608-2623.
- Han X, Zhong Y, Cao L, Zhang L (2017) Pre-trained alexnet architecture with pyramid pooling and supervision for high spatial resolution remote sensing image scene classification. Remote Sensing 9: 848.
- Latif S, Qadir J, Farooq S, Imran MA (2017) How 5G wireless (and concomitant technologies) will revolutionize healthcare? Future Internet 9: 93.
- Saari P (2015) Music mood annotation using semantic computing and machine learning (Doctoral dissertation, University of Jyväskylä).
- Mondada L (2016) Challenges of multimodality: Language and the body in social interaction. Journal of sociolinguistics 20: 336-366.
- Anadon LD, Chan G, Harley AG, Matus K, Moon S, et al. (2016) Making technological innovation work for sustainable development. Proceedings of the National Academy of Sciences 113: 9682-9690.
- Nyati S (2018) Revolutionizing LTL Carrier Operations: A Comprehensive Analysis of an Algorithm-Driven Pickup and Delivery Dispatching Solution. International Journal of Science and Research (IJSR) 7: 1659-1666.
- Demirkan H, Spohrer J (2014) Developing a framework to improve virtual shopping in digital malls with intelligent self-service systems. Journal of Retailing and Consumer Services 21: 860-868.
- Chojceki P (2020) Artificial Intelligence Business: How you can profit from AI. Przemek Chojceki.
- Karl M (2018) Risk and uncertainty in travel decision-making: Tourist and destination perspective. Journal of Travel Research 57: 129-146.
- Hollebeck, LD, Macky K (2019) Digital content marketing's role in fostering consumer engagement, trust, and value: Framework, fundamental propositions, and implications. Journal of interactive marketing 45: 27-41.
- Gokhale T, Banerjee P, Baral C, Yang Y (2020) Vqa-lol: Visual question answering under the lens of logic. In European conference on computer vision Cham: Springer International Publishing 379-396.
- Dicianno BE, Joseph J, Eckstein S, Zigler CK, Quinby E, et al. (2018) The voice of the consumer: a survey of veterans and other users of assistive technology. Military medicine 183: e518-e525.
- Poria S, Hazarika D, Majumder N, Mihalcea R (2020) Beneath the tip of the iceberg: Current challenges and new directions in sentiment analysis research. IEEE transactions on affective computing 14: 108-132.
- Ramakrishnan S, Agrawal A, Lee S (2018) Overcoming language priors in visual question answering with adversarial regularization. Advances in Neural Information Processing Systems 31.
- Das A, Anjum S, Gurari D (2019) Dataset bias: A case study for visual question answering. Proceedings of the Association for Information Science and Technology 56: 58-67.
- Ławrynowicz A (2020) Creative AI: A new avenue for the Semantic Web?. Semantic Web 11: 69-78.
- Sampat SK, Yang Y, Baral C (2020) Visuo-linguistic question answering (VLQA) challenge. arXiv preprint arXiv: 2005.00330.
- Gill A (2018) Developing A Real-Time Electronic Funds Transfer System for Credit Unions. International Journal of Advanced Research in Engineering and Technology (IJARET) 9: 162-184.
- Pochangou PM (2020) Internet of Things: technologies and applications in healthcare management and manufacturing (Doctoral dissertation, Université Laval).
- Kafle K, Kanan C (2017) Visual question answering: Datasets, algorithms, and future challenges. Computer Vision and Image Understanding 163: 3-20.
- Nyati S (2018) Transforming Telematics in Fleet Management: Innovations in Asset Tracking, Efficiency, and Communication. International Journal of Science and Research (IJSR) 7: 1804-1810.

Copyright: ©2022 Vedant Singh. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.