# EXPLOITING SEMANTIC KNOWLEDGE FOR IMAGE CAPTIONING USING DEEP LEARNING

By

**Ali Raza**



**NATIONAL UNIVERSITY OF MODERN LANGUAGES**

**ISLAMABAD**

**2024**

# Exploiting Semantic Knowledge for Image Captioning Using Deep Learning

By

**Ali Raza**

MSSE, National University of Modern Languages, Islamabad, 2024

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF

THE REQUIREMENTS FOR THE DEGREE OF

**MASTER OF SCIENCE**

In **Software Engineering**

To

FACULTY OF ENGINEERING AND COMPUTING



NATIONAL UNIVERSITY OF MODERN LANGUAGES ISLAMABAD

Ali Raza,2024

# THESIS AND DEFENSE APPROVAL FORM

**The undersigned certify that they have read the following thesis, examined the defense, are satisfied with overall exam performance, and recommend the thesis to the Faculty of Engineering and Computing for acceptance.**

**Thesis Title:** EXPLOITING SEMANTIC KNOWLEDGE FOR IMAGE CAPTIONING USING DEEP LEARNING

**Submitted by:** Ali Raza          **Registration #:** 48 MS/SE/F21

Master of Science in Software Engineering

Degree name in full

Software Engineering

Name of Discipline

Dr. Jaweria Kanwal          _____

Name of Research Supervisor          Signature of Supervisor

Dr. Sumaira Nazir          _____

Name of HOD (SE)          Signature of HOD (SE)

Dr. Noman Malik          _____

Name of Dean (FEC)          Signature of Dean

1st April, 2024

Date

# AUTHOR'S DECLARATION

I <u>Ali Raza</u>

Son of <u>Akhtar Hussain</u>

Registration # <u>48 MS/SE/F21</u>

Discipline <u>Software Engineering</u>

Candidate of **Master of Science in Software Engineering (MSSE)** at the National University of Modern Languages do hereby declare that the thesis **Exploiting the Semantic Knowledge for Image Captioning Using Deep Learning** submitted by me in partial fulfillment of MSSE degree is my original work, and has not been submitted or published earlier. I also solemnly declare that it shall not, in the future, be submitted by me for obtaining any other degree from this or any other university or institution. I also understand that if evidence of plagiarism is found in my thesis/dissertation at any stage, even after the award of a degree, the work may be canceled and the degree revoked.

_____

Signature of Candidate

<u>  Ali Raza            </u>

Name of Candidate

<u>1<sup>st</sup> April, 2024        </u>

Date

# AUTHOR'S DECLARATION

I <u>Ali Raza</u>

Son of <u>Akhtar Hussain</u>

Registration # <u>48 MS/SE/F21</u>

Discipline <u>Software Engineering</u>

Candidate of **Master of Science in Software Engineering (MSSE)** at the National University of Modern Languages do hereby declare that the thesis **Exploiting the Semantic Knowledge for Image Captioning Using Deep Learning** submitted by me in partial fulfillment of MSSE degree is my original work, and has not been submitted or published earlier. I also solemnly declare that it shall not, in the future, be submitted by me for obtaining any other degree from this or any other university or institution. I also understand that if evidence of plagiarism is found in my thesis/dissertation at any stage, even after the award of a degree, the work may be canceled and the degree revoked.

_____

Signature of Candidate

<u>  Ali Raza            </u>

Name of Candidate

<u>1st April, 2024        </u>

Date

# ACKNOWLEDGEMENT

# DEDICATION

ALHAMDULILLAH… All gratitude be to Allah Almighty for molding me into the person I am today and allowing me to realize my ambition. I dedicate this thesis work to my supervisor, ***"Dr. Jaweria Kanwal"*** for guiding me and giving me her precious time whenever I needed it the most during my Research Study. Your encouragement, understanding, and sacrifices throughout my academic journey have shaped me into the researcher I am today. Thank you for always being there for me.

To my dear parents Your unconditional love, support, and guidance have been the bedrock of my academic achievements. Through your sacrifices and endless encouragement, you have shown me what true dedication and hard work look like. I am forever grateful for the sacrifices you have made to provide me with an education, and for instilling in me a passion for learning. It is your unwavering commitment to excellence that has inspired me to strive for greatness in all that I do.

# ABSTRACT

The technique of generating textual explanations for images is commonly referred to as image captioning. It has attracted a lot of attention recently because it may be used in a variety of fields. There are some challenges in image captioning, one of them is the lack of incorporating semantic knowledge in generating image captions. Semantic knowledge can be helpful in object detection by exploiting relationships among objects and in language semantics. In this study, the issue of image captioning is investigated by combining two efficient models, the vision transformer (ViT) and the generative pre-trained transformer 2 (GPT-2). The ViT uses self-attention techniques that are applied to image patches to capture visual elements and overall context from images. The GPT-2 model complements ViT with extraordinary language production abilities that enable it to produce content that is cohesive and related to the situation. An encoder-decoder-based deep learning model is proposed where the ViT performs the encoder function, extracting meaningful visual representations from images, while the GPT-2 model performs the decoder function, producing descriptive captions based on the retrieved visual features. This method makes it possible to seamlessly combine textual and visual information, producing captions that faithfully reflect the content of the input images. The potential of this combination is demonstrated through empirical analyses, highlighting the advantages of utilizing both language and visual components in the 'image captioning' process. My research strengthens multimodal AI systems by bridging the gap between visual and language comprehension. The experiments were performed on the MS COCO dataset and Flicker 30k dataset. The model was validated using various evaluation metrics. Results show an improvement as Bleu-1, Bleu-2, Bleu-3, Bleu-4, Rogue, and Meteor by 10.58, 20.45, 21.07, 34.19, 0.3, and 11.16 respectively. The other evaluation metrics like Meteor improved by 11.16 and the Rogue metric improved by 0.3 on the MS COCO dataset.

# Table of Contents

# List of Figures

# List of Tables

# CHAPTER 1

# INTRODUCTION

The realms of computer vision and natural language processing are increasingly intrigued by the captivating and intricate task of generating automatic descriptions for images. Humans can identify the objects in an image and their spatial relationships only by looking at them. Humans then develop a natural language description of that image. Although it is simple for people, there are various phases when it comes to implementing it for machines. The deep learning techniques process for the generation of image descriptions is known as 'Image Captioning'.

The computer vision community is actively researching the problem of automatically creating image captions [1]. It combines two of the key fields of artificial intelligence: computer vision and natural language processing [2]. An image caption generation model must not only be able to identify the objects within a picture, but also be able to articulate their relationships in natural languages [3], [4]. The creation of image captions is a difficult undertaking since it requires figuring out the existence and relationships of various items as well as organizing human-like sentences to represent this information [5]

An image consists of a lot of information that can be perceived differently by various people from a different perspective [2], [4]. Image processing is crucial for content-based image retrieval (CBIR), which has numerous applications in digital libraries, online searching, biomedical, business, the military, and education [4]. Social media sites like Facebook and Twitter have the use to create descriptions automatically from images. Image details created from social media sites can include the location (such as a beach or cafe), and what we are wearing. and most importantly, what we are doing? It is being widely used in many applications, including the text descriptions retrieved from image searches or image look-ups of the sentences provided [3], [6].

Generating descriptions for photographs through automated processes is known as image caption creation [7]. It entails comprehending the image's semantics, which calls for knowledge of the primary objects, their varied characteristics, attitudes, and interactions with the image. To provide appropriate captions, it must also deduce the underlying semantic meanings. Figure 1.1 displays several images accompanied by their respective captions. The captions "A couple of kids walking around with colorful umbrellas", "A green bird standing on peeled bananas in a background", and "A man in a soccer uniform playing soccer on a field" The provided captions correspond to the visuals presented in Figures 1.1(a), (b), and (c), respectively [8].

Generating descriptions for images holds significance for a variety of purposes. This technology, for instance, can be used to construct picture search engines, intelligent computer-human interactions, and automatic image captioning. Platforms like Facebook and Twitter may automatically produce descriptions based on a photograph, our location (such as a beach or cafe), what we are wearing, and what we are doing. Additionally, it can be applied to summarize events. Figure 1.2 provides some instances of how captioning has been used in various contexts, including (a) scene description for people who are blind, (b) interaction among humans and robots, and (c) text-based picture retrieval.

Captioning for images is a crucial field of study. Both visual comprehension and a language description for that image are necessary for automatic caption generation. One of the major issues in computer vision is image understanding. Natural Language Understanding (NLU) includes language description [9]. An image encoder that extracts

features from an image and a language decoder that creates captions for that picture make up a standard 'Image Captioning' architecture [8].



(a) A couple of kids walking around with colourful umbrellas.

(b) A green bird standing on peeled bananas in a background.

(c) A man in a soccer uniform playing soccer on a field.

**Figure 1.1** Examples of appropriate Generated Caption

## 1.1   Image understanding

An image needs to be described properly to understand its meaning. It is difficult to automatically describe an image's content in natural language. It must accurately describe the object and their relationships to prevent incorrect information. Figure 1.1 describes the basic concept of image captioning. The description of the image should cover all the aspects of an image to prevent missing any important information [2]. There are various sources of images, including television, websites, social media sites, and news channels. It is very difficult to generate a caption automatically [10][11]. An image description should explain not just the elements of the image but also their relationships to one another, as well as their characteristics and the activities they are engaged in. For creating a description for an image, it is necessary to ensure the interconnection among the objects, actions, and things in an image.

A human can understand images without any technical assistance because the human being has the cognitive abilities that help in paying attention to the image, perceiving the image according to the language, and recalling from his memory. After seeing an image, a human makes a perception about the image in his mind according to his previous knowledge and makes a thought about that image. However, the whole process is difficult

for the machine to understand the image. For this computer vision and natural language processing techniques are used. To understand these concepts, an example was followed that humans are waiting at the station for the train, but the train is not at the station. A human can understand by seeing the background location but the machine cannot detect the exact location without detectable aspects [2], [10].

## 1.2   Natural Language Understanding

As stated by the natural language understanding (NLU) perspective, text creation involves several steps. To begin, the elements of the input must be comprehended, often known as content selection, the material must be organized, which is referred to as text planning; and finally, must verbalized, which is referred to as surface realization. Surface realization encompasses lexicalization, which includes the selection of appropriate words, production of referential expressions utilizing correct pronouns, and subsequently, aggregation, a process that merges pertinent information [12]. To extract meaning, context, and intent from text or voice, a computer or AI system must first be able to read and interpret human language. This process is known as natural language comprehension. This is a crucial element of natural language processing (NLP), which enables computers to properly engage with and respond to human language.

## 1.3   Human Cognition

The mental processes and skills that allow people to learn, process, store, and apply knowledge are referred to as human cognition. It includes several cognitive functions, including perception, attention, memory, language, problem-solving, reasoning, and decision-making. Human cognition is impacted by variables like heredity, environment, and experience and involves intricate relationships between various cognitive processes. To better understand how people perceive, comprehend, and solve issues, cognitive psychology, and cognitive science investigate human cognition. The basic processes of human cognition are being studied by researchers using experimental approaches, brain imaging tools, and computational models [13] [14].

**Figure 1.2** Steps of the Human Cognition Process

## 1.3.1 Steps of Human Cognition

There are following steps included in the cognition process are given below:

**Attention:** The ability to focus on a specific external stimulus is a cognitive function. The ability to concentrate on specific environmental information while disregarding other details is referred to as attention. As attention has limitations in both capacity and duration, it is essential to efficiently oversee the available attention resources to comprehend the surroundings effectively [13] [14].

**Language:** Comprehending and conveying concepts using spoken and written language are a pair of cognitive activities integral to linguistic growth and progression. This enables interaction with other people and is crucial for thought [13] [14].

**Learning:** Cognitive functions involving the assimilation of novel information, its synthesis, and its amalgamation with pre-existing knowledge are imperative for the process of learning. Learning is an experience-based alteration of behavior that often lasts a long time. It involves acquiring fresh knowledge and developing novel skills. When contemplating the act of learning, it's often directed towards the structured schooling received by children and young adults. But learning happens continuously throughout life and is not only limited to the classroom [13] [14].

**Memory:** Memory, a fundamental cognitive process, facilitates the encoding, storage, and retrieval of information. This crucial function is integral to the learning process and empowers individuals to maintain knowledge about both the external environment and their personal histories. The capacity of human memory includes the capacity to store and retrieve data. Sometimes data is simply not stored in memory correctly and people misremember or forget things [13] [14].

**Perception:** The process of interpreting sensory data from the environment, such as visual, auditory, and tactile stimuli, to develop accurate representations of the outside world is referred to as perception. People use perception, a cognitive process, to gather information from their senses and use it to react to and engage with their environment. Our sensory perception of the world is referred to as perception. It is the process of becoming aware of things and connections through the use of our senses. A deeper understanding of our environment was gained through such experiences. Our capacity to perceive our surroundings is contingent upon the mental mechanisms that utilize to decipher data, such as utilizing memory to identify the visage of a friend or detecting the scent of a familiar perfume [13] [14].

**Thought:** Every cognitive process depends on thought. It enables people to make decisions, solve problems, and use higher-level reasoning. Humans evaluate, create ideas, generate ideas, and build concepts through mental processes known as thoughts. They are essential to human cognition because they allow people to properly analyze information, address issues, make decisions, and communicate [13] [14].

## 1.4  Machine Cognition

The ability of artificial intelligence systems to understand, observe, and reason about the world in a way that is similar to human cognition is referred to as machine cognition in deep learning. It entails creating algorithms and models that are capable of higher-level cognitive functions including perception, comprehension, reasoning, and decision-making in addition to more conventional pattern recognition and prediction tasks. This comprises activities like text summarization, question-answering, sentiment analysis, and semantic comprehension [15], [16]. By enabling machines to learn from massive volumes of data and automatically extract complicated characteristics and representations, deep learning plays a significant role in machine cognition. Multiple-layered deep neural networks may capture hierarchical and abstract representations, enabling machines to comprehend and interpret data in a more sophisticated way [17][18].

**Machine Perception:**   The ability of machines to collect data from a variety of sources, including detectors, cameras, microphones, and other data inputs, is referred to as perception. Like human senses gather sensory data, this stage enables machines to collect unprocessed data from the environment**.** The process of perception entails creating models that can recognize and interpret sensory input, incorporating areas such as speech recognition, natural language processing, and visual perception (object recognition, scene understanding), and other types of sensory data analysis [18].

**Feature Extraction:** Machines analyze the raw information and extract pertinent features or patterns. This process entails locating crucial components and traits that aid machines in better comprehending the data. Identifying visual features in images or linguistic features in natural language processing are two examples of feature extraction in different fields [18].

**Learning and representation:** Following the extraction of the features, machine learning algorithms are used. These algorithms learn from the data and build representations of it using the retrieved features. The computer may use supervised learning (with labeled data), unsupervised learning (without labels), or reinforcement learning, depending on the job, to enhance its comprehension of the data [18].

**Reasoning and inference:** The process of reasoning entails applying the acquired representations to reach deductive and inductive conclusions. Logic, probabilistic thinking, and inductive or deductive reasoning techniques are frequently used in this procedure. Machines can examine data patterns and draw conclusions or predictions based on previously learned information [18].

**Problem Solving and Decision Making:** Machine cognition entails creating models that can decide wisely and take appropriate action based on the facts at hand. This encompasses activities like autonomous navigation, reinforcement learning, and decision-making in challenging contexts. The goal of machine cognition in deep learning is to close the cognitive gap between artificial intelligence systems and those possessed by humans.

Machines can perform decision-making and problem-solving activities using their acquired knowledge and reasoning skills. Based on the input and their internal comprehension, they may assess many options, balance benefits and drawbacks, and make the best decisions [18].

'Image Captioning' is the field of computer vision and NLP that are used for the generation of textual descriptions of an image. This technology is an important part of this era and is used in various industries and domains like social media, the medical field, especially for impaired people, advertising and publication of news, etc. As this technology is used in various industries and domains, there needs to be more accuracy in this area so that users can use this technology properly. Different researchers work in this domain but there is still a lack of accuracy and performance in this field. Reduction in the accuracy and performance is due to the lack of semantic knowledge means that the relationships among the objects within the image and categorization of the regions of the objects according to their relationships.

(a) Scene description for visually impaired people.

(b) Human-robot interaction.

**Figure 1.3** Image Captioning Applications

According to the problem of this domain 'Image Captioning'. An encoder-decoder framework was proposed based on the two-transformer model Vision Transformer and Generative Pre-trained Transformer 2. These two transformer models are combined into a ViT-GPT-2 encoder-decoder framework that mainly focuses on the semantic knowledge of the image to improve the accuracy and performance of the generated image caption. The patching procedure was changed from 16x16 pixels to 8x8 pixels because the number of patches increased and feature extraction and fetching semantic knowledge from each patch of the image from each patch of the image. Semantic knowledge of the image and feature extraction improves and the accuracy and performance of the image caption also improves by this procedure of ViT-GPT-2 encoder-decoder framework

## 1.5   Problem Statement

Currently, there is a significant focus on the research of image captioning. There are many studies conducted on image captioning. Many researchers use various methods for image captioning but natural language descriptions that are generated are not accurate. The researcher did not achieve the level of accuracy in the caption generation [19], [20]. The

lack of accuracy in the image caption is due to the lack of semantic knowledge of an image. Semantic knowledge needs to improve the relationship among the objects within the image and focus on the regions of the object and their relationship with each other [19], [20].

## 1.6 Research Objectives

The research objectives are given below:

Obj1: To improve the performance of caption generation for an image through deep learning.

Obj2: To exploit semantic knowledge for image captioning.

## 1.7 Research Questions

In order to fulfill these research aims; the following inquiries have been taken into account:

RQ1: How can a deep learning model be proposed to generate a more accurate caption of an image?

RQ2: How can the semantic knowledge of an image be exploited in the process of image captioning?

## 1.8 Scope of Study

Research on image captioning, which entails analyzing images and providing a verbal description for them, is a rising field. Understanding an image involves more than just identifying and locating items; it also entails knowing the sort of scene or environment, the placement of the objects, and how they relate to one another. Syntactic and semantic knowledge of the language is necessary for creating well-formed sentences. This is for the understanding of images for disabled persons who cannot feel or understand the scenes. This research helps blind persons.

## 1.9   Research Contribution

The primary objective of this research is to generate an improved caption of an image to describe the accurate content of the image. The following are major contributions to achieving the research objectives and goals.

**Integrated Approach:** A transformer-based new approach is introduced that uses two different transformer models such as Vision Transformer (ViT) and Generative Pre-trained 2 (GPT-2), allowing for a combination of visual and textual information.

**State-of-the-art Findings:** After thorough experimentation on benchmark datasets, including MS COCO and Flicker 30k. ViT-GPT-2 fusion consistently demonstrates superior performance when compared to current leading models. Enhancements are evident in quantitative measures like BLEU, METEOR, CIDER, and ROUGE scores, as well as in qualitative evaluations of the generated captions. The performance of this approach improves as can see from the results that come from the experiment. Qualitative and quantitative analysis of the results are given in Chapter 5.

## 1.10 Summary

In this chapter, a complete overview of 'image captioning' is discussed. The basic understanding of the image and natural language and the concepts of the cognitive process were also discussed, and their steps for humans as well as the machine to mimic the image understanding like a human being. In this chapter, the research problem was further discussed, including my research objectives and research question my proposed solution, and my contribution to the conduction of research in the domain of 'Image Captioning'.

# Chapter 2

# BACKGROUND

## 2.1 Machine learning

A field within the realm of artificial intelligence (AI) known as machine learning gives systems the ability to learn from concepts and information without having to be explicitly programmed. To anticipate the characteristics and trends present within data, enhancing future results and informed choices commences with initial observations, including firsthand experiences. Deep learning relies on a set of machine learning methods. that use many nonlinear transformations to represent high-level abstractions in the data [21].



**Figure 2.1** Process Steps for Machine Learning

## 2.2   Deep learning

A popular and exciting area of machine learning is deep learning. Deep learning has gained popularity in the research of machine learning because of its accuracy, and performance The most efficient machine learning method in terms of performance, supervision, cost, and time is deep learning. Deep learning transcends being a constrained learning approach; instead, it encompasses various approaches and structures that prove valuable in addressing a diverse array of intricate problems. Operating in a highly intricate manner, this methodology acquires distinctive and discriminative features. The evolution of deep learning methodologies has been noteworthy, demonstrating impressive performance across a spectrum of applications enhanced by robust security protocols [21]. Leveraging the backpropagation technique, it stands out as the optimal choice for discerning complex architectures within high-dimensional data. Deep learning is widely used in business, science, and government because it has made great strides and performs admirably in a wide range of applications, this encompasses adaptive testing, classification of biological images, computer vision, detecting cancer, processing natural language, recognizing objects and faces, deciphering handwriting, identifying speech, analyzing the stock market, shaping smart cities, as well as various other applications [21].

| Input | → | Feature Extraction + Classification | → | Out Put |

**Figure 2.2** Process Steps for Deep Learning

## 2.3   Convocational Neural Network (CNN)

Convolutional Neural Networks (CNNs) represent a category of neural networks that have demonstrated remarkable performance in tasks such as image recognition and categorization. Their exceptional accuracy in image classification establishes convolutional neural networks as the predominant approach in the field of computer vision.

**Figure 2.3** Convolutional Neural Network (CNN) Model [22]

Convolutional Neural Networks (ConvNets), leverage input images to integrate distinct features into the structure of the network. As a result, this leads to a decrease in the quantity of parameters and an improvement in the execution efficiency of the forward function implementation. ConvNets consist of neurons that possess biases and weights. These neurons receive inputs and perform a dot product computation. The structure of convolutional neural networks sets them apart from regular networks. In conventional neural networks, inputs traverse through multiple layers of concealed neurons. Each layer comprises a cluster of neurons, and these neurons within each layer establish complete connectivity with all neurons in the subsequent layer. Furthermore, the neurons in each layer function autonomously without interconnections. The ultimate layer, known as the output layer, signifies the network's predictions and maintains complete connectivity as well [22], [23].

**Figure 2.4** Convolution Neural Network and Simple Neural Network [23]

In Figure 2.4, the illustration on the left displays a common three-layer neural network, while the depiction on the right showcases a Convolutional Neural Network (CNN) organizing its neurons using three dimensions: height, width, and depth [22].

The convolution operation is a mathematical procedure carried out within a convolutional neural network. It encompasses a mathematical computation where two functions (referred to as 'a' and 'b') are combined to generate a third function. This operation is dependent 'a*b' and is computed by integrating the product of the two functions, where one of them is mirrored and shifted.

$$(a * b)(t) = \int_{-\infty}^{+\infty} f(\tau)b(t - x)d\tau \qquad (2.1)$$

The process of convolution encompasses three key components [22], [23]:

- The initial image, which serves as the input to the operation.
- A feature detector, often referred to as a "kernel" or "filter," designed to identify distinct features within the image. This is typically represented as a matrix with dimensions like 5x5 or 7x7.
- The resulting feature map, also denoted as an activation map, illustrates the locations of particular features within the image. The term "feature map" is used due to its function in mapping out the positions of specific features within the image.

Convolution Mathematics: The convolution's output, influenced by the input factor, is elucidated as stated in reference [52].

Imagine that:

Let's consider an image with an input size of N, denoting both its width and height as N, resulting in an N * N image size. If there is a filter of size F * F, a stride value of S, and zero padding of P, the resulting image size can be represented as O.

$$O = \frac{N-F+2P}{S} + 1 \tag{2.2}$$

CNN Model's Fundamental Component: The essential structure of a CNN consists of a series of layers, each utilizing a unique function to transform a set of activations into a different form. CNN designs are formulated through three primary categories of layers.

## 2.3.1 Convolutional Layer

An essential aspect of a Convolutional Neural Network (CNN) architecture is the convolutional layer. This element conducts convolutions on input images by employing a set of filters or kernels. These filters can detect specific characteristics in the image, like edges, textures, and patterns. Individually, each filter convolves with the input image, producing a feature map that indicates the existence of the particular attribute that the filter aims to identify. Aggregating the results of multiple filters produces a multi-channel feature map. Later on, these feature maps undergo one or more non-linear activation functions, such as ReLU, to add a degree of non-linearity to the model. The outcome is then directed to the next layer, which might consist of another convolutional layer or a pooling layer. This series of actions is repeated in an iterative manner until the ultimate output is attained.

**Figure 2.5** Convolutional Operation of CNN [23]

The characteristic enhances computational efficiency by reducing the overall parameter count and simplifying the computation process.

The control of the output volume of the convolution layer is determined by three parameters as listed below:

**Depth:** In the initial layer of a CNN model, the dimensionality of the input volume signifies the number of color channels found in the input image. For a colored image, this dimensionality equals 3, corresponding to the red, green, and blue channels. In scenarios involving grayscale or black and white images, the dimensionality becomes 1. Meanwhile, the dimensionality of the resultant volume is established by the quantity of filters employed on the input image.

**Stride:** Stride is employed to traverse the width and height of the given image. The filters shift one pixel at a time when the stride is set to 1. Alternatively, with a stride of 2, the filters advance by 2 pixels for each movement

**Zero Padding:** Zero-padding involves adding zeros to the input image within the input layer. This approach helps control the dimensions of the input layer. Without implementing zero-padding, there is a risk of losing certain edge-related features [22], [23].

## 2.3.2  Pooling Layer

Pooling layers are employed to diminish the dimensions of feature maps by condensing the characteristics within a specific region. This action subsequently results in a reduction of the requisite learnable parameters. The pooling layer condenses the features detected within a segment of the feature map produced by a convolutional layer [23].

**Pooling Function:** Pooling is executed by employing diverse methods to decrease the input parameters. These techniques encompass computing the average, minimum, or maximum values solely within the sub-regions.

Varieties of Pooling Methods: The following list enumerates various types of pooling functions.

- **Max Pooling:** This procedure yields the maximum attainable value.
- **Average Pooling:** It computes the average and furnishes the maximum value.
- **Weighted Average Pooling:** By considering the pixel's proximity to the center, it calculates the weight of the surrounding area.
- **L2 Norm Pooling:** The outcome is the square root of the region encompassing the neighborhood's rectangles.

Most ConvNet designs incorporate Max Pooling as a technique to decrease computational expenses [23].

## 2.3.3  Fully Connected Layer

Much like a neural network, each neuron within the fully connected layer establishes connections with every other neuron situated in the underlying layer. The determination of its activation involves matrix multiplication incorporating both weight and bias factors [23]. The final levels in the CNN architecture, these layers are in charge of making the final predictions. Each neuron in these fully connected neural networks is linked to every other neuron in the layer below it. A SoftMax layer receives the output of the final fully connected layer and uses it to classify data into many categories. Neurons within the fully

connected layer establish direct connections with neurons in the two contiguous layers, without forming links with any neurons within those particular layers. This arrangement bears similarity to the positioning of neurons in conventional models of Artificial Neural Networks (ANNs) [22], [23].



**Figure 2.6** Fully Connected Layer of CNN [22]

## 2.4  Activation Functions

An Activation Function is responsible for determining the activation status of a neuron. This entails utilizing fundamental mathematical operations to assess the significance of the neuron's input to the network. Some of the activation functions that are commonly used are given below:

## 2.4.1  ReLU (Rectified Linear Unit)

The abbreviation for Rectified Linear Unit is ReLU. It adjusts the threshold of the input value to zero, yielding the input number for positive values and zero for negatives. Upon incorporation as an activation function, ReLU exhibited a speed increase of sixfold

compared to the function in the Alex Net architecture. The subsequent equation represents the ReLU formula [23].

$$f(x) = max(0, x)$$



**Figure 2.7** ReLU Function's Curve [23]

The advantages of employing the Rectified Linear Unit (ReLU) function comprise the following:

- Its efficient computation facilitates swift network convergence.
- Despite its seemingly linear shape, it introduces non-linearity.
- The feasibility of its incorporation into backpropagation stems from the presence of its derivative function.

## 2.4.2 SoftMax

A generalized version of the logistic function for multiple classes is known as the SoftMax activation function. This function takes an input vector and produces a probability distribution that encompasses all available classes. This becomes particularly valuable when dealing with multi-class classification challenges, where the objective is to categorize an input into one among numerous potential classes. In the context of neural networks with multiple layers intended for classification purposes, the SoftMax function is commonly employed in the output layer (described as Equation 2.3). By utilizing this function, the probabilities for each class are computed, and the sum of these probabilities across the

entire output layer totals 1. In essence, the SoftMax function aids in predicting the class to which the input most likely pertains, rather than making a discrete selection between classes [23].

$$SoftMax(x) = \frac{e^{x_k}}{\sum_{i=1}^{m} e^{x_i}} \tag{2.3}$$

This method is frequently employed as the last layer in an image classification convolutional neural network (CNN). The outcome of this ultimate layer represents a probability distribution across all categories, and the network's prediction is determined by selecting the class with the utmost probability [23].

A Convolutional Neural Network (CNN) is comprised of multiple stages of convolutional and pooling operations. These are succeeded by one or multiple fully connected layers, culminating in the inclusion of a SoftMax layer. The convolutional and pooling stages work harmoniously to extract distinctive attributes from the initial image, while the fully connected layers amalgamate these characteristics to construct a succinct yet meaningful depiction of the image. Ultimately, the SoftMax layer generates a distribution of probabilities across all categories, thereby facilitating predictions [22], [23].

### 2.4.3 Sigmoid

Artificial neural networks, including convolutional neural networks, commonly employ the sigmoid activation function. This particular type of squashing function was utilized to transform input values into a range spanning from 0 to 1. This characteristic rendered it valuable in tasks involving binary outcomes, like categorizing input into one of two classes, or in probability modeling.

$$S(A) = \frac{1}{1+e^{-a}} \tag{2.4}$$

Here, x represents the input value provided to the activation function. The sigmoid function exhibits a gradual shift from 0 to 1, facilitating improved convergence of the neural network while it undergoes the training phase.

**Figure 2.8** Sigmoid Function's Curve [23].

When employing a CNN for binary classification, it is common to employ the sigmoid function as the activation mechanism for the output layer [23].

## 2.5 Recurrent Neural Network (RNN)

Recurrent Neural Networks (RNNs) represent a type of artificial neural network (ANN) wherein neurons are enabled to establish cyclic connections among themselves as well as with other neurons within the identical layer. The structure depicted in Figure 2.9 illustrates a fundamental RNN configuration. Distinct examples of RNN architectures encompass the Simple Recursive Networks (SRNs) and the Long Short-Term Memories (LSTMs). The utilization of cyclical connections proves advantageous in grasping the step-by-step anticipation of results. In this scenario, the present outcome hinges on not only the current input but also the preceding outcomes, alongside the current input. In contrast to conventional FNNs, RNNs exhibit greater effectiveness in encapsulating such Markov models with their enhanced elegance [24].

**Figure 2.9** An RNN exhibiting the distinctive cyclic interconnections [24]

## 2.5.1  Simple Recurrent Network (SRN)

An uncomplicated recurrent network featuring cyclical connections within its layers is referred to as a simple recurrent network (SRN). The diagrams presented in Figures 22a and 22b depict the structural makeup of an SRN. These visuals portray identical architectures, with the latter offering a deeper understanding of the functioning of an RNN. It does so by elucidating the processes taking place during each time step and illustrating how the former hidden layer's output influences both the current hidden output and the ongoing input. In contrast, the former diagram demonstrates the conventional representation, highlighting the recursive connection. Notably, the output from the previous time step exerts influence on the current output owing to its reliance on the former hidden state(s) [24], [25].

**Figure 2.10** Architecture of SRN [24]

a) An SRN architecture with clearly named inputs, outputs, and weight factors. (b) Using an unrolled SRN, visualize the effect of earlier concealed states on current output.

**Forward Propagation in SRN:** Let's consider an activation function denoted as $f$. In the context of a Sequential Recurrent Network (SRN), the variables $x_t$, $h_t$, $h_{t-1}$, and $y_t$ which represent the current input, current hidden state, previous hidden state, and current output value, respectively. The architecture consists of three weight matrices: $Wxh$ for input-to-hidden, $W_{hh}$ for hidden-to-hidden, and $W_{hy}$ for hidden-to-output transitions. Additionally, denote $in_{y0}$ and $in_{h0}$ as the input values to the activation functions in the output and hidden layers respectively. With these components established, the forward pass output in an SRN can be represented using the following sequence of equations [24].

$$in_{h0} = (W_{xh}x_t + W_{hh}h_{t-1}) \tag{2.5}$$

$$h_t = f(in_{h0}) \tag{2.6}$$

$$in_{y0} = W_{hy}h_t \tag{2.7}$$

$$in_{y0} = W_{hy}h_t \tag{2.8}$$

Assuming initiation at time step 't' and continuation until 't+3', the SRN demonstrates progression. It takes inputs $xt, xt+1$, $xt+2, xt+3$ and yields outputs $yt, yt+1$, $yt+2, yt+3$.

The process involves intermediary hidden node results $h_t$, $h_{t+1}$, $h_{t+2}$, $h_{t+3}$, depicted in Figure 2.11. The computation of these values (from equations 2.9 to 2.17) leads to this illustration. During the first-time step, the previous hidden output is treated as having a value of zero. The same weights are utilized in all subsequent iterations.



**Figure 2.11** Illustration of Forward propagation [24]

**Backpropagation in SRN:** Aside from following the instructions provided for backpropagation in Feedforward Neural Networks (FNNs) as detailed in section 3.2, the backpropagation procedure includes an extra time-related element. This is frequently denoted as backpropagation through time (BPTT). The underlying concept is that in contrast to a hidden node within an FNN, a hidden node within a Recurrent Neural Network (RNN) at any given time step bears direct responsibility for the incurred cost not only at that specific time step but also at the subsequent time step. Furthermore, it also indirectly contributes to any subsequent errors that arise [24]. Thus, drawing from the principles of BPTT and leveraging equations (2.9), (2.10), (2.11), and (2.12), The BPTT equations for the scenario depicted can be formulated in Figure 2.11.

$$\frac{\partial c}{\partial h_{t+3}} = \frac{\partial c3}{\partial h_{t+3}} = W_{hy} \cdot df(i_{y3}) \cdot \frac{\partial c3}{\partial y_{t+3}} \tag{2.9}$$

$$\frac{\partial c2}{\partial h_{t+2}} = W_{hy} \cdot df(in_{y2}) \cdot \frac{\partial c2}{\partial y_{t+2}} \tag{2.10}$$

$$\frac{\partial c1}{\partial h_{t+1}} = W_{hy} \cdot df(in_{y3}) \cdot \frac{\partial c3}{\partial y_{t+3}} \qquad (2.11)$$

$$\frac{\partial c0}{\partial h_t} = W_{hy} \cdot df(in_{y3}) \cdot \frac{\partial c3}{\partial y_{t+3}} \qquad (2.12)$$

$$\frac{\partial c}{\partial h_{t+2}} = \frac{\partial c2}{\partial h_{t+2}} + \frac{\partial c3}{\partial h_{t+2}} = \frac{\partial c2}{\partial h_{t+2}} + W_{hh} \cdot df(in_{h3}) \cdot \frac{\partial c}{\partial h_{t+3}} \qquad (2.13)$$

$$\frac{\partial c}{\partial h_{t+1}} = \frac{\partial c1}{\partial h_{t+1}} + \frac{\partial c2}{\partial h_{t+1}} = \frac{\partial c1}{\partial h_{t+1}} + W_{hh} \cdot df(in_{h2}) \cdot \frac{\partial c}{\partial h_{t+2}} \qquad (2.14)$$

$$\frac{\partial c}{\partial h_t} = \frac{\partial c0}{\partial h_t} + \frac{\partial c1}{\partial h_t} = \frac{\partial c2}{\partial h_t} + W_{hh} \cdot df(in_{h1}) \cdot \frac{\partial c}{\partial h_{t+1}} \qquad (2.15)$$

$$\frac{\partial c}{\partial W_{hh}} = h_t \cdot df(in_{h1}) \cdot \left( \frac{\partial c1}{\partial h_{t+1}} + \frac{\partial c2}{\partial h_{t+1}} \right) \qquad (2.16)$$

$$\frac{\partial c}{\partial W_{xh}} = x_t \cdot df(in_{h0}) \cdot \left( \frac{\partial c0}{\partial h_{t+1}} + \frac{\partial c1}{\partial h_{t+1}} \right) \qquad (2.17)$$

The term $df(x)$ represents the differentiation of the activation function f, computed at the specific value of x. Additionally, c0, c1, c2, and c3 symbolize the expenses accrued during time intervals ranging from t to t+3 [24].

**Vanishing and exploding gradients:** The challenge of the vanishing gradient arises from how gradients are retrogressed over time within SRNs. This leads to swift multiplication of the derivative term, $d_f$, associated with the activation function as depicted in equations (2.13) through (2.14) being repeatedly substituted into (2.15). While the provided instance employs a time step of four to enhance understanding, it's important to note that the number of time steps in RNNs is commonly well beyond four. Unfortunately, in cases where the derivative of the activation function is less than one, the gradients diminish rapidly over time. On the contrary, when the derivative exceeds one, the gradients experience exponential growth. The probability of experiencing the challenge of vanishing gradients becomes more pronounced as time elapses. This implies that when a substantial temporal gap exists between the dependencies of the output and preceding inputs, the gradients will attenuate at an earlier stage. Consequently, the model faces difficulties in accurately capturing temporal dependencies. This underlying concern led to the development of RNNs, also known as Recurrent Neural Networks. Weight initialization is commonly carried out by setting weights to have an average of zero and a standard deviation of 0.001.

Moreover, as widely used activation functions like the logistic sigmoid and hyperbolic tangent typically yield derivatives that remain at or below one, the occurrence of excessive gradient growth is infrequent. Conversely, the primary obstacle revolves around the vanishing gradient phenomenon. This challenge is effectively addressed when employing activation functions such as ReLU and piecewise linear activations, given that their derivatives consistently hold values of either one or zero. This property effectively mitigates the vanishing gradient predicament. Hence, the utilization of these activation functions helps alleviate the problem. The challenge of vanishing gradients also emerges in deep FNNs, underscoring why CNNs exhibit improved performance with RELU activations. To address the potential occurrence of excessively large gradient values, practitioners commonly apply gradient clipping, which limits them from surpassing a specific threshold [24], [25].

**Inability to capture long-term dependencies:** SRNs suffer from an incapability to grasp extended connections between data points, attributed to challenges like the vanishing gradient issue detailed earlier. Additionally, the tendency for more recent input values to dominate and replace prior hidden states within the network contributes to a diminished influence of these preceding states on the overall learning process. This phenomenon is depicted in Figure 2.12. The outcomes of the video classifier prioritize the present input frame, as evidenced by its tendency to revise predictions when the current input frame significantly deviates from previous ones. Consequently, this leads to the neglect of the cumulative effects of prior output predictions [24].

LSTMs tackle this concern by integrating input gates that manage the impact of the previous hidden state and the current input state on the present hidden state. Correspondingly, output gates are utilized to govern the effect of the current cell state on the immediate output, as illustrated in Figure 2.13. The input gates, shown in the color blue, function to avoid the dominating influence of the current input, ensuring that the importance of prior hidden states is not excessively undermined. Simultaneously, the output gates, distinguished in red, regulate the influence of the cell state on the output by deferring predictions about the activity until the conclusion of video frames [24], [26].

**Figure 2.12** Demonstration of the challenge in capturing extended dependencies through a video-input-based activity classifier employing an SRN approach [24].



**Figure 2.13** Maintaining Extended Relationships through Input and Output Gates [24]

## 2.5.2 Long Short-Term Memory (LSTM)

As a substitute for SRN, Long Short-Term Memories (LSTMs) often exhibit superior performance in capturing temporal dependencies due to their multiple gating mechanisms. LSTMs have exhibited their higher performance compared to RNNs in different tasks within the field of natural language processing. These tasks encompass handwriting recognition (Source: [20]), language translation (Source: [22]), as well as image and video annotation [24]. The conventional LSTM cell comprises three essential gates: the input gate, the forget gate, and the output gate. Both a standard RNN (SRN) cell and a simplified LSTM cell layout are illustrated in Figure 26. Instead of traditional neurons, an LSTM-based architecture employs these LSTM cells to construct the hidden layers. The subsequent discussion will elaborate on the distinct functions fulfilled by the different gates within LSTM units [24].



**Figure 2.14** Internal Architecture of LSTM [24]

By moderating the impact of the present input and the output from the previous time step on the current cell state, the input gate grants the LSTM the capability to uphold or overwrite the information from the prior hidden layer and the current input node. Even though LSTMs demonstrate proficiency in accommodating temporal dependencies across more than 1000-time steps, there are occasions when it becomes essential for the network to eliminate past input to avert the incorporation of undesirable dependencies during the learning process. The forget gate provides a method to diminish or conceivably completely disregard the influence that a former cell state exerted on the current cell state. The output gate manages the movement of output from the present cell state to the output of the current hidden state [24], [26].

**Forward Propagation Problem:** The updates for LSTM at every time step $x_t$ can be outlined as per references [26]. This is in the context of a sequential input series $[x1, x2,...$ $xt-1, xt,..., xT]$, wherein $x_t-1$ and $x_t$ represent consecutive inputs in the system.

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \tag{2.18}$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \tag{2.19}$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \tag{2.20}$$

$$g_t = \phi(W_{xg}x_t + W_{hg}h_{t-1} + b_g) \tag{2.21}$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t \tag{2.22}$$

$$h_t = o_t \odot \phi(c_t) \tag{2.23}$$

In this particular context, $g_t$ serves as the input node or modulation gate responsible for modifying the input received by the LSTM cell. It mirrors the outcome of the input gate. The results of the forget gate and output gate are denoted as $f_t$ and $o_t$, respectively. The memory cell, labeled as $c_t$, emerges from a combination of the previous memory cell influenced by the forget gate and the input node managed by the input gate. Lastly, the current hidden state known as $h_t$ represents the activation of the current cell state, which is under regulation. Operating as controllers that govern the flow of signals, all gates utilize the sigmoid activation function to produce outputs like $i_t$, $f_t$, $o_t$, and so forth. The utilization of the tanh activation function is a widespread approach within Long Short-Term Memory

(LSTM) networks. Using backpropagation through time, the iterative adjustment of weights (W) and bias values (b) becomes achievable. An LSTM network has the potential to incorporate multiple LSTM layers, where each layer commonly consists of numerous cells (nodes), mirroring the concept of artificial neural networks (ANNs) with their numerous concealed layers, each comprising a multitude of hidden nodes [24], [26].

**Avoiding the vanishing gradient problem:** Sustaining a consistent error loop enables the mitigation of the vanishing gradient issue in Long Short-Term Memory networks (LSTMs) [5]. The gradient propagated backward at the previous time step, t-1, is ascertained using formulas (46), given the gradient backpropagated to the cell state at time t, denoted as ctt.

$$\frac{\partial c_{t-1}}{\partial (t-1)} = f(t) \odot \frac{\partial c_t}{\partial t} \tag{2.24}$$



**Figure 2.15** Depiction showcasing the BPTT process applied to an LSTM cell [24].

The gradient continues to exist as the output of the forget gate remains close to 1. When the forget gate equals 0, any connection between previous time steps and the current ones disappears. This makes the vanishing gradient acceptable. For a comprehensive depiction of the computation, please consult Figure 2.11.

## 2.6  Attention

In deep learning, the term "attention" refers to a method that enables models to concentrate on particular input data segments while processing information. It helps the model to choose focus on pertinent data by allowing it to assign various levels of priority to various input items. Different neural network topologies now include attention mechanisms as a critical component, which helps them perform better on tasks like machine translation, 'Image Captioning', question answering, and more.

The attention mechanism computes attention weights or scores for each location or region using the feature map or feature vectors as input. Each region's value or relevance in the visual input is indicated by its score. In many neural network topologies, notably in the areas of natural language processing (NLP) and computer vision, attention mechanisms play a vital role. They enable models to concentrate on various input data elements with varied degrees of relevance. The fundamental concept behind attention is to give various weights to various input data components to represent their relative importance. In NLP, for instance, attention processes assist the model in determining which words in the source phrase to concentrate on while creating each word in the target sentence [27], [28].

**Context Vector:** A context vector that depicts the attended regions or features is created by computing a weighted sum of the feature vectors using the attention weights. A concentrated representation of the input image is given by the context vector. Different designs and techniques, such as soft attention or hard attention, can be used to implement visual attention. Soft attention employs continuous attention weights, which enables the model to pay attention to several regions at once. Hard attention, on the other hand, employs discrete attention, picking a specific area or patch at each stage.

## 2.6.1  Visual Attention

Visual attention is a strategy that deep learning models employ to analyze information while selectively focusing on particular areas of an input image. It enables the model to highlight important portions of the image while ignoring less important ones. In

tasks like image production, object detection, and image captioning, visual attention processes have been particularly effective.

Visual attention, commonly referred to as spatial attention, describes a model's or system's capacity to focus only on particular areas or aspects of visual input. It draws inspiration from how people see images, where paying attention is key to focusing on important details in a scene. To enhance performance and interpretability, visual attention mechanisms have been effectively applied to a variety of computer vision applications, including picture categorization, object identification, and image captioning [29].

When generating judgments, neural networks can concentrate on particular areas of an input image thanks to a mechanism called visual attention that was inspired by human perception. To enhance their performance on tasks like 'Image Captioning', object identification, and image classification, numerous deep learning models, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have widely embraced this method [29], [30]:

## 2.6.2  Semantic Attention

Semantic attention, often referred to as content-based attention or soft attention, is a technique frequently employed in sequence-to-sequence models for things like text generation, image captioning, and machine translation. The model can concentrate on various input sequence segments while creating the output sequence thanks to this attention technique. It aids the model in evaluating the significance of various input sequence components for each stage of output creation [28]. When creating captions for images, a technique known as semantic attention focuses on the semantic significance of certain areas or items within the image. Semantic attention, in contrast to conventional attention processes, considers both the semantic meaning and relationships between things to create more accurate and contextually appropriate captions. To understand the semantic connections between objects and how they affect caption creation, semantic attention models are used. To direct the attention mechanism, they make use of additional semantic information such as verbal signals or labels for object categories [16], [31].

## 2.7 Transformer

In the article Attention Is All You Need, Transformers were presented [32]. It serves as the foundation for some popular versions like GPT-2 [33] and BERT [34]. Language translation models and question-and-answer-based models are two examples of transformer models' many applications. Transformers can be utilized for a variety of application cases because of their versatile architecture [35].

Figure 2.16 illustrates the structure of a Transformer model. The image showcases the Encoder stack with N identical layers on the left side and the Decoder stack with N identical layers on the right side. To gain a comprehensive understanding of the roles of each side of the transformer, it's essential to first familiarize ourselves with the objectives of the individual subtasks outlined within the tiers.

### 2.7.1 Layers of Transformer

**Multi-head attention layer:** Take the phrase "The animal didn't want to cross the street because it was tired" as an example. Transformer employs an "Attention" mechanism to help a machine link "it" to the animal. A Transformer can focus on other aspects from input that are closely connected to the feature it is now focusing on thanks to attention, which establishes a relationship between the two. Using embedding methods, words in NLP must be converted into vectors. The vector for a sentence of length n would be n x d_model, where d_model is the feature encoding size. The initial stage of the encoder/decoder receives these vectors and applies a sequence of "Scaled-Dot Product Attention," also known as "Multi Head Attention," to them.

**Figure 2.16** Transformer Architecture

To summarize the function of attention, it can be said that one input information is allowed to be concentrated on by a transformer while assessing how other features are closely related to it. establishing a connection between them. Scaled-Dot Product Attention is implemented in the movie Transformers.

To summarise the function of attention, It can be said that one input information is allowed to be concentrated on by a transformer while assessing how other features are closely related to it. establishing a connection between them. Transformers uses "Scaled-Dot Product Attention" as its application of attention. Q, K, and V, or query, key, and value vectors, are built for each word in the sentence, which can be characterized as:

- Q: The vector awaiting determination
- K: The vector symbolizing the features
- V: The vector representing the true input values

Figure 2.16 illustrates the realization of the 'Scaled-Dot Product Attention' mechanism.



**Figure 2.17** Scale Dot Product Attention

According to the information presented in Figure 4.2, the process begins with the computation of the dot product between the Q and K vectors. Subsequently, this dot product is divided by the square root of $dk$, which serves the purpose of addressing the issue of potential gradient explosion. The resulting step involves the application of a SoftMax function to the normalized dot product, thereby deriving the weights necessary for the scaling of V. The summarized expression for the Scaled-Dot Product Attention can be expressed using equation 2.25.

$$\text{Attention}\ (Q, K, V) = \text{softmax}\ \left(\frac{QK^T}{\sqrt{d_k}}\right) V \tag{2.25}$$

A Multi-Head Attention (MHA) layer comprises a set of 'h' Self-Attention heads, with each head independently computing the Scaled Dot Product Attention. The outcomes of these computations are then merged through concatenation, resulting in the Multi-Head Attention value. This procedure is illustrated in the diagram provided.

## Multi-Head Attention



**Figure 2.18** Multihead Attention can be described as a parallelized version of Scaled Dot Product Attention.

Utilizing multiple attention heads enhances the effectiveness of the attention layer. This enables the model to concentrate on various positions, providing multiple sets of weight matrices for the Q, K, and V components. The attention computed at each step is sometimes referred to as 'z-score' or simply 'z.' In a scenario where, for instance, 8 attention heads are employed, the resulting values would be denoted as ($z0$, $z1$, $z2$, $z3$, $z4$, $z5$, $z6$, $z7$). However, these values cannot be directly fed into the subsequent layer of the encoder/decoder. A transformation is necessary, involving the concatenation of all outcomes, followed by multiplication with the additional weight matrix WO. It's worth noting that WO is a supplementary weight matrix that undergoes joint training with the model.

**Figure 2.19** The procedure for acquiring the Z matrix, which then serves as the input for the subsequent Transformer stages.

The mathematical definition of Multi Head Attention can be expressed as follows.

$$\text{MultiHead}\ (Q, K, V) = \ \text{Concat}\ (\text{head}\ _1,\ \text{head}\ _2, \ldots,\ \text{head}\ _h) W^o \qquad (2.26)$$

$$\text{where head}\ _i = \ \text{Attention}\ \left(QW_i^Q, KW_i^K, VW_i^V\right) \qquad (2.27)$$

**Add and Norm Layer:** The layer executes two activities, as its name suggests. The 'Add' portion of the process, which controls flow via residual connections, is the initial phase. 'Norm', the next step, accomplishes layer normalization. Hence, the result of this layer will align with the subsequent equation.

$$\text{Add}\ \backslash\&\ \text{Norm}\ = \ \text{LayerNorm}\ (x +\ \text{Sublayer}\ (x)) \qquad (2.28)$$

Here, x represents the input to any sublayer (either MHA or Feed Forward), and Sublayer(x) denotes the resultant output.

**Feed Forward Layer:** Each layer has a fully connected point-wise feed-forward network. It uses ReLU activation to carry out two linear transformations. This layer determines the weights used during exercise. This particular stratum computes the weights in the course of training. It can be formally expressed in mathematical terms as stated below:

$$FF(x) = \text{ReLU}\,(xW_1 + b_1)W_2 + b_2 \qquad (2.29)$$

$$ReLU\,(x) = \max(0, x) \qquad (2.30)$$

In this context, W1 and W2 represent matrices of network weights, while b1 and b2 indicate bias terms.

## 2.7.2  Positional Encoding

The embeddings' relative or absolute positions are introduced into the model by the Transformers using positional encoding. This preserves the parallel execution format of the token sequence. The sine and cosine representations of the position and training parameters are used to calculate the values for positional encoding. To construct position-aware embeddings, the positional encodings are combined with the language features.

$$PE_{(pos,2i)} = \sin\left(\,pos\,/1000^{2i/d\_\text{model}}\right) \qquad (2.31)$$

$$PE_{(pos,2i+1)} = \cos\left(\,pos\,/1000^{2i/d\_\text{model}}\right) \qquad (2.32)$$

Here, 'i' represents the embedding's position, and 'd_model' indicates the embedding's dimension.

## 2.7.3  Transformer Encoder-Decoder

**Encoder Transformer:** The encoder is made up of N repetitions of the following layers: Multi-Head Attention, Add & Norm, Feed-Forward, and Final Add & Norm. Positional encoding is added to the input of the encoder layer, which is in the form of n x d_model, where n is the total number of characteristics and d_model is the size of the features' embedding. One layer of the encoder stack receives this input only.

**Figure 2.20** Structure of the Encoder: The Encoder utilizes Self Attention where Q, K, and V are set to be the same.

**Decoder Transformer:** The transformer's decoder portion has N identical layers as well. The Decoder is made up of three sublayers: the 'Mask Self Attention layer', the 'multi-head attention' between the "encoder output" and the output of the Masked Self Attention layer, and the 'Feed Forward Layer'. An Add & Norm layer comes after each sublayer. Because it establishes a connection between encoder output and decoder input, the

intermediate MHA layer of the decoder is also known as the encoder-decoder attention layer.



**Figure 2.21** Internal Composition and Arrangement of the Decoder Stack. The Decoder comprises two distinct variants of Multihead Attention layers.

The initial decoder layer employs Masked Self-attention, which restricts Self-Attention to focus solely on preceding positions. The input to the decoder is right-shifted and subjected to masking. This arrangement ensures that the word at the ith position can only attend to preceding words, up to position i. This masking strategy, facilitating parallel computations, safeguards against the decoder gaining advanced knowledge and engaging in premature optimization during training—a form of 'cheating'. The accompanying illustration depicts the application of masking, effectively concealing future embeddings from the Self-Attention mechanism. To execute masked self-attention, a lookahead mask is introduced, taking the form of a lower triangular matrix.

**Figure 2.22** Mask Self-Attention's Working

The subsequent stratum constitutes a conventional Multi-Head Attention layer. Nonetheless, it employs the Encoder's output as the 'K' and 'V' matrices, while utilizing the output from the Masked Self Attention layer as matrix 'Q'. This results in the nomenclature 'Encoder-Decoder MHA layer'. This stratum enables full interconnections between all decoder positions and input sequence positions. The ultimate tier within this arrangement is the Feed Forward network.

## 2.7.4  Linear & SoftMax Layer

Similar to previous seq2seq models, this layer is also. The output of the decoder is transmitted through a fully connected linear layer, which projects a size n x vocab_size output, where n is the expected sentence length and vocab size is the size of the language's vocabulary. The SoftMax layer is then applied to the resultant matrix, giving each word in the output sentence a probability distribution across the lexicon.

## 2.8  Object Detection

An essential objective within the field of computer vision is object detection, which revolves around the identification and spatial localization of items within an image that are of significance. Deep learning techniques have made great strides in the field of object

detection recently, enabling precise and effective recognition across a variety of item categories [36].

Commonly, one- or two-stage approaches are used by deep learning-based object detection algorithms:

**Two-stage object detection:** Two-stage detectors first provide several area proposals—potential item bounding boxes—before classifying and improving them. The classification step assigns object labels and fine-tunes the bounding box coordinates, while the region proposal stage chooses a collection of possible regions based on object ness scores or other factors. Support vector machines (SVMs) are used for classification in R-CNN, which uses a CNN to extract area information. These efforts prepared the ground for further advancements like Fast R-CNN and Faster R-CNN [36].

**One-stage object detection:** Through the direct prediction of item bounding boxes and class probabilities, one-stage detectors carry out object detection in a single step. Compared to two-stage techniques, these detectors are often faster but may not be as accurate. To localize objects, they employ anchor boxes, predetermined bounding box forms at different scales, and aspect ratios. In the realm of single-stage object detectors, the lineage of models known as You Only Look Once (YOLO) stands out as a noteworthy architectural advancement. YOLO fashions. Due to its quickness and ease of use, this real-time detection method has grown in popularity [37].

## 2.9   Feature Extraction

In deep learning, the process of obtaining useful representations or features from the initial raw input data is referred to as feature extraction. It entails taking informative visual information from pictures or video frames for use in computer vision processes including object identification, image classification, and image segmentation. Because of their ability to construct hierarchical and distinguishing representations, deep learning models, specifically convolutional neural networks (CNNs), have been widely employed for feature extraction. [38].

The process of feature extraction involves the subsequent actions:

**Data Preprocessing:** To ensure that raw input data is suitable for feature extraction, preprocessing is frequently necessary. Data cleansing, normalization, dimensionality reduction, or addressing missing values are some examples of the tasks that may be involved in this step. Preprocessing seeks to improve the data's quality and get rid of any noise or extraneous information that can make it difficult to extract features.

**Feature Selection:** The process of feature selection entails selecting a subset of pertinent features from the original data. This process is essential because it lowers the data's dimensionality by removing duplicate or pointless information that could cause overfitting or increase computing complexity. The selection of features can be done using a variety of ways, including statistical methods, correlation analysis, or domain expertise.

**Feature Encoding:** Following feature extraction, the features are frequently encoded into a suitable numerical format that machine learning algorithms can understand. To ensure interoperability with the selected learning method, this stage entails translating the features into a standardized representation, such as vectors or matrices. One-hot encoding, bag-of-words representation, and vector embeddings are a few common encoding strategies.

## 2.10 Vision Transformer (ViT)

The Transformer paradigm [39], which was initially created for natural language processing, is now used for computer vision problems via the Vision Transformer (ViT) architecture. When it comes to image recognition tasks, ViT obtains cutting-edge performance and, in some cases, outperforms conventional convolutional neural networks (CNNs).

### 2.10.1 Image Patching

The input image is divided into fixed-size, non-overlapping patches as the first stage of ViT. The Transformer model's input is a vector representation that has each patch linearly embedded into it. These patches' sizes are hyperparameters that can be altered

according to the particulars of the task. A typical patch size, for instance, is 16x16 pixels or 8x8.

## 2.10.2 Patch Embeddings

Learnable positional embeddings are added to the patches that are linearly embedded. The model can comprehend the relative positions of various patches in the image thanks to these positional embeddings, which encode spatial information. The input information for the following Transformer encoder layers is created from a combination of patch embeddings and positional embeddings.



**Figure 2.23** Working Pattern of Vision Transformer

### 2.10.3 Positional Embeddings

The learnable positional embeddings are added to the patch embeddings to incorporate spatial information and the relative placements of the patches. The model can comprehend the geometric relationships between patches thanks to these positional embeddings. Patch and positional embeddings are combined to create the input data for the Transformer encoder layers.

### 2.10.4 Transformer Encode Layer

Similar to what is done in jobs involving natural language processing, ViT makes use of a stack of Transformer encoder layers. A feed-forward neural network and a multi-head self-attention mechanism make up the two sub-layers of each encoder layer.

- **Multi-Head Self-Attention**: This technique enables each patch in the image to pay attention to neighboring patches and record cross-patches interdependence. To enable the model to concentrate on pertinent patches while processing the data from each patch, it computes attention weights between all pairs of patches.

- **Feed-Forward Neural Network:** The patch embeddings travel via a feed-forward neural network after the self-attention sub-layer, which performs non-linear modifications on the representation of each patch.

- **Layer Normalization and Residual Connections**: Each sub-layer (self-attention and feed-forward) is followed by layer normalization and residual connections, much like the original Transformer architecture, to stabilize and assist the training process.

- **Classification Head:** The classification head is the last layer in the ViT model, and it uses the output from the last Transformer encoder layer to convert it to the number of classes needed for the particular task. To obtain class probabilities, the classification head for image classification typically consists of a fully connected layer followed by a SoftMax activation function.

## 2.10.5 Pre-training of Vision Transformer

ViT is trained using a self-supervised objective on a sizable dataset during the pre-training phase. Given the context of the other patches in a picture, the model learns to anticipate corrupted or masked patches inside the image. ViT can learn visual representations that capture crucial elements in photos thanks to this technique.

## 2.10.6 Fine- tunning of Vision Transformer

ViT may be fine-tuned for a variety of downstream tasks, including segmentation, object detection, and image classification. The classification head is trained on a dataset specific to the task with labeled samples during fine-tuning, adapting the model's parameters to the intended job.

## 2.11 Generative Pretrained Transformer -2 (GPT-2)

Modern language model GPT-2 (Generative Pre-trained Transformer 2) [40] was created by OpenAI. It is a Transformer architecture version made specifically for jobs involving natural language processing. GPT-2 is a massive language model with 1.5 billion parameters that can produce text that is coherent and appropriate to the situation.

## 2.11.1 The architecture of GPT-2

- The Transformer design, which comprises a stack of encoder-decoder layers, is the foundation of GPT-2.
- GPT-2 uses only the decoder portion of the transformer, in contrast to conventional sequence-to-sequence transformers. It is an autoregressive language model, which means that each token that is formed is dependent upon the prior tokens that were generated.

**Figure 2.24** Working Pattern of Generative Pre-trained Transformer

## 2.11.2 Transformer Architecture of GPT-2

GPT-2 is built on the Transformer design, which was first presented in the paper "Attention Is All You Need." [28] It is a neural network architecture that significantly depends on mechanisms for maintaining self-attention to effectively handle data sequences. The Transformer's essential parts include the following:

- **Multi-Head Self-Attention**: With the use of this attention mechanism, the model can process each location while concentrating on various portions of the input sequence. The importance of various tokens in the context of the entire sequence is captured by computing attention weights for each input token concerning all other tokens.

- **Feed-Forward Neural Networks:** The model processes the attention outputs from the attention layers using feed-forward neural networks to build higher-level representations.

- **Layer Normalization and Residual Connections:** Each sub-layer (attention and feed-forward) is followed by layer normalization and residual connections to aid training and handle the vanishing/exploding gradient problem.

### 2.11.3 GPT-2 Specifics

GPT-2 lacks the encoder component found in devices like BERT because it is a decoder-only Transformer. It processes input sequences using a stack of Transformer decoder layers. A language modeling objective is used to pre-train the model on a huge corpus of text data in an unsupervised manner. Predicting the next word in a sentence based on the context of the previous words is a part of the pre-training process. The pre-training exercise aids the model's acquisition of contextual representations that capture linguistic patterns in both syntactic and semantic terms.

### 2.11.4 Tokenization

- The text is tokenized into smaller units, such as words or sub-words, before being fed into the model (using, for instance, byte-pair encoding).
- The embedding layer is then used to transform each token into a fixed-size vector representation.

### 2.11.5 Pre-training of GPT-2

- Using an unsupervised learning strategy, the GPT-2 has been pre-trained on a sizable corpus of text data.
- The model learns to anticipate the possibility of the following word in a sentence given the context of the preceding words during pre-training.
- Given the previous context, the training objective is to increase the likelihood of the target term. The model gains knowledge of grammar, syntax, and semantics from the large amount of text material it encounters thanks to this process.

## 2.11.6 Fine-tuning of GPT-2

Following pre-training, GPT-2 can be improved for a variety of downstream tasks. For fine-tuning, the model is trained on a smaller dataset relevant to the desired job, such as sentiment analysis, text production, or translation. GPT-2 can be adjusted to more accurately carry out particular tasks by making minor adjustments.

## 2.11.7 Self-Attention Mechanism

- At the core of the Transformer architecture lies the self-attention mechanism, which allows the model to evaluate the importance of different words in context with each other.
- By establishing long-range dependencies and capturing global context, the model is better able to recognize the connections between the words in a phrase.

## 2.11.8 Autoregressive Generation

- The model generates new text one token at a time during inference or text creation, starting with a seed text (sometimes referred to as a prompt).
- The model uses the feed-forward layers and self-attention mechanism of the decoder at each step to forecast the likelihood of the upcoming token on the information provided by the tokens generated earlier.
- GPT-2 can generate text using a variety of sampling techniques, including temperature-based sampling, beam search, and greedy sampling. With greedy sampling, the token with the highest probability is selected at each step. With beam search, the top k likely sequences are tracked.
- The diversity and coherence of the resulting text are impacted by different sampling procedures.

The key benefits of GPT-2 include its capacity for producing text that is coherent and contextually suitable, as well as its language comprehension abilities and adaptability to different NLP tasks via fine-tuning.

## 2.11.9 Loss Function

GPT-2, like many other autoregressive language models, uses the cross-entropy loss function during training. In the context of language modeling, the cross-entropy loss measures the dissimilarity between the model's predicted probability distribution over the vocabulary and the actual distribution of the next token in the sequence [27], [28].

Mathematically, given a sequence of input tokens $x_1, x_2, \ldots, x_n$ and their corresponding target tokens $y_1, y_2, \ldots, y_n$ and sample strategies where $y_i$ is the next token in the sequence after $x_i$ the cross-entropy loss $L$ is computed as:

$$L = -\frac{1}{n}\sum_{i=1}^{n}\log p(y_i \mid x_1, x_2, \ldots, x_{i-1})$$

Where $p(y_n|x_1, |x_2, \ldots, x_{i-1})$ represents the probability assigned by the model to the target token $y_i$ given the proceeding input tokens $x_1, x_2, \ldots, x_{i-1}$.

During training, the goal is to minimize this cross-entropy loss, which encourages the model to assign high probabilities to the correct tokens in the sequence and penalizes it for making incorrect predictions [27], [28].

## 2.12 Summary

In this chapter, background about different concepts (terminologies and models) that were used in previous studies literature review like CNN, RNN, and transformer models like Vision Transformer and Generative Pre-trained Transformer 2 are discussed. Some terminologies like machine learning, deep learning, semantic attention, visual attention,

and feature extraction were also discussed. The reader who reads this document can understand these terminologies and models that are used in the research.

# CHAPTER 3

# Literature Review

## 3.1  Introduction

In this section, the literature review concerning research in our specific field of image captioning was explored. A literature review is an examination and synthesis of previously published research and academic articles that are pertinent to a given topic. Examining a literature review's content, structure, and conclusions is part of the analysis process. This research literature review is about 'Image Captioning'. The distinction between the information conveyed by images and words presents a semantic disparity, despite the abundant explicit and implicit visual semantic details often present within images. It is challenging to accurately express an image's visual information using only human language.

In general, two paradigms may be used to classify existing systems for creating image captions: retrieval-based and generation-based. By collecting comparable photos from the training dataset per similarity metrics, retrieval-based image captioning techniques generate a caption for an input query image. The extracted feature vectors are compared one to another in a typical calculation. The best candidate image's caption is then applied to the input image [5].

Deep learning-based techniques can manage this challenging task, according to recent work on 'Image Captioning' [20], [30], [41]. The primary idea of these methods,

which are typically based on the encoder-decoder framework from machine translation [20], [42], [43], is to approach the task of creating image captions as translating an image into a text description.

## 3.2 Deep Learning Techniques for Image Captioning

A semantic depiction known as the scene graph, rooted in graphs, establishes a connection between photographs and their corresponding natural descriptions. They can't see the entire horizon at once because they are captioning. Instead, they move their focus from one place to another in a fluid, sequential manner. Typically, their attention is drawn to larger, more colorful objects first, then to smaller ones. Their strategy entails creating a scene-graph-oriented representation that mimics the human attention mechanism. The process of image captioning was divided into two independent stages known as idea cognition and sentence creation to include the scene graph as a middle step [6]. In their method, they build a set of semantic ideas and then introduce the CNN-RNN-SVM architecture to produce a sequence based on scene graphs. After that, this sequence is transformed into a bit vector and used as the RNN's next phase's input [6].

A combination of an image encoder and a language decoder constitutes this architectural configuration. This innovative design was first introduced by works [20], [44] in their work. They adeptly adapted the encoder-decoder architecture originally used for machine translation to align to generate image captions. This involved utilizing a CNN to serve as the dedicated image encoder. Further enhanced the decoding capabilities by replacing the standard RNN with the long short-term memory (LSTM) architecture, resulting in significant advancements in generating image captions. Since then, numerous researchers have built upon this framework, focusing on enhancing both components to improve image captioning  [20], [44].

This article's main contribution is an improved visual attention model. At first, they suggest computing the focus intensity coefficient of the attention mechanism using the model's context data at each timestep. They then utilize this coefficient to automatically alter the focus intensity to more precisely retrieve visual data. They also include topic terms

related to the image scene to describe the semantic information of the scene, which is then included in the language model. The attention mechanism is employed to identify the visual and semantic information about the scene that the model concentrates on at each timestep, combining it to produce captions that are more accurate and relevant to the scene. [20].

This article proposes a brand-new method for captioning images called Visual-Semantic Double Attention (VSDA). A unique Semantic Attention (SEA) technique was employed to extract semantic features, and their approach comprises two primary components: initially, an adapted visual attention model is employed to capture image information from sub-regions. The specific significance of each attribute word is frequently ignored by traditional attribute-based models, which instead combine them into recurrent neural networks. As a result, there are too many unnecessary semantic traits present. Simply put, VSDA's fundamental power resides in its ability to make good use of semantic features while also minimizing the impact of unnecessary attribute words, hence improving the accuracy of semantic guidance[45].

This paper [46] examines the active learning technique used in image captioning and introduces a novel model called the Structural Semantic Adversarial Active Learning (SSAAL) model. SSAAL utilizes visual and textual data from both local and global viewpoints to identify the most informative samples The model comprises three components: a semantic constructor responsible for generating a structural representation of an image, an SC supervisor that oversees the representation using a multi-task learning approach, and a labeled / unlabelled state discriminator that utilizes adversarial learning to differentiate between samples with varying descriptions.

In [47] their primary objective is to connect the realms of vision and language by enhancing image features with textual concepts. This approach establishes a strong foundation for describing images. They investigated the use of textual representations for image features to describe prominent regions within images at a textual level. Their method incorporates the Textual Distilling Module and Textual Association Module, which aim to harness vast and enhanced textual data to enhance the comprehension of images. Through extensive experiments conducted on the well-known Flickr30k and MSCOCO image captioning datasets, they successfully validated the efficacy of their approach.

In the publication [48], a novel image captioning approach is presented, referred to as the domain-specific image captioning generator. This innovative model employs a combination of visual and semantic attention mechanisms to craft descriptive captions for images. The paper refers to the output as the "general caption," which represents a caption generated for a given image. Moreover, the model can produce domain-specific captions by substituting specific terms within the general caption with domain-specific vocabulary, leveraging semantic ontology principles. The model's performance was assessed through comprehensive qualitative and quantitative experiments. Nevertheless, a limitation of their approach is identified in its inability to seamlessly integrate semantic ontology in an end-to-end manner.

In [49] They suggest an alternative to the current image-caption pair replacement approach, which suffers from the issue of inaccurate pseudo-labeling. The proposed solution is called the Image-Caption Pair Replacement Algorithm (I-CPRA) and consists of two sub-modules: the bounding box scaling algorithm and the two-stage semantic graph structure. More specifically, within the image replacement strategy, they aim to overcome the challenge of strict resolution and aspect ratio replacement requirements between novel objects and source objects. To tackle this, they introduced a bounding box scaling algorithm as part of our approach. The number of potential objects for the model's novel object replacement is expanded, increasing synthetic images. To tackle the issue of contextual logical errors caused by the direct replacement of descriptions, their solution involves a two-stage semantic graph structure. This structure aims to minimize phrase collocation errors within the context by relying on associations of co-occurring semantic adjacency.

They presented a novel approach for generating image captions using an adversarial training technique. To enhance the precision of the captions, they incorporated a semantic filter module to extract valuable contextual information from semantic vectors. Their model utilizes a two-separated LSTM architecture to acquire both global and local image features along with semantic vectors. By employing adversarial training, the resulting captions seamlessly integrate precise details and are expressed smoothly and fluently. Their model aims to acquire meaningful semantic knowledge by considering the contextual information within image regions. This approach leads to enhanced fluency and accuracy

in the generated captions. In terms of fluency, they observed an average reduction of 14.1% in the PPL score compared to the baseline model Semantic Filter - Generative Adversarial Networks (SF-GAN) [50].

The objective of this study [51] is to present an innovative model that makes use of the Fully Convolutional Network – Long Short-Term Memory (FCN-LSTM) framework for generating a finely detailed attention map at the grid level. The new model guarantees that the main object solely influences the visual characteristics of each grid cell. Through the integration of labels specific to each grid (semantic segmentation), the visual features of diverse grid cells are interconnected. This strategy empowers the model to aptly capture wider contextual details from semantic labels, thus encompassing supplementary semantic context. As a result, this technique enhances the contextual information accessible to the language LSTM decoder, furnishing a more comprehensive comprehension of the input. By creating a mechanism of precise and contextually guided visual attention, a connection is established between relevant visual information and the corresponding semantic meanings in the text. Through three experiments that encompass qualitative and quantitative analyses, our model demonstrates the ability to produce high-quality captions. Specifically, it achieves remarkable levels of accuracy, comprehensiveness, and diversity.

In the research conducted by the study [19], a new and innovative deep architecture named the attention-based Encoder-Decoder is presented. This architectural design makes effective use of convolutional features derived from a CNN model that underwent pre-training on ImageNet. These features are combined with object features that were extracted from the pre-trained You Look Only Once (YOLO) v4 model, trained on MS COCO. Furthermore, the paper introduces a novel method of positional encoding referred to as the "importance factor," which is applied to object features. The performance of the proposed model was assessed using both the MS COCO and Flickr30k datasets and a comparative analysis was conducted against similar approaches. The incorporation of this original technique for feature extraction led to a notable increase of 15.04% in the CIDEr score.

This publication [52] presents an innovative algorithm Attention Model - Fully Connected Network (ATT-FCN) that combines two methodologies using a semantic attention model. The algorithm adeptly concentrates on pertinent semantic concept

proposals, integrating them into the concealed states and yields recurrent neural networks. This merging of selection and integration creates an iterative loop linking the higher-level and lower-level computations. To gauge its effectiveness, they executed assessments on two extensively utilized benchmarks: Microsoft COCO and Flickr30K. The outcomes of our experiments consistently illustrate that our algorithm outperforms existing cutting-edge techniques across a range of evaluation criteria.

The present study [3] introduces a clear-cut and easily comprehensible reasoning framework Visual Semantic Reasoning Network (VSRN) for generating visual representations that encompass significant objects and semantic concepts present within a scene. The strategy involves establishing associations among diverse regions of an image while employing Graph Convolutional Networks to generate characteristics that encompass semantic connections. Furthermore, the authors recommend the integration of gate and memory mechanisms to globally process semantic reasoning on these relationship-amplified characteristics. This facilitates the identification of distinguishing information and the gradual formation of a comprehensive scene representation. Through experimentation, the effectiveness of the approach is confirmed in achieving a fresh pinnacle of performance in image-text matching on the MS-COCO [28] and Flickr30K [40] datasets. The findings underscore its superiority over the current leading approach, displaying relative enhancements of 6.8% for image retrieval and 4.8% for caption retrieval on MS-COCO (utilizing the Recall@1 [3] benchmark on a 1K test grouping). Furthermore, on Flickr30K, the proposed model attains relative improvements of 12.6% for image retrieval and 5.8% for caption retrieval (Recall@1).

In the work cited as [53], scene graph-base captions (SG2Caps) are presented as a framework that leverages only scene graph labels to achieve competitive performance in the task of image captioning. The central goal is to minimize the difference in meaning between two scene graphs: one obtained from the input image and another generated from the associated caption. This is executed by integrating information about the spatial positioning of objects and labels related to human interaction (HOI), leading to the creation of an additional HOI graph. SG2Caps surpasses existing models that rely solely on scene graphs for captioning, highlighting the potential of scene graphs as a promising strategy

for enhancing image captioning. This method directly uses scene graph labels, avoiding the computationally intensive graph convolutions on high-dimensional CNN features, resulting in a significant reduction of 49% in trainable parameters.

This study [54] presents an innovative structure for generating image captions by integrating features with enriched semantics and introducing difficult counterexamples for improved performance. The method suggested in this research merges these components through a Semantic-Enhanced Module, which consists of a sub-network for matching images and text, along with a Feature Fusion layer. This combination produces semantically enriched features that are imbued with extensive semantic details. Additionally, the authors propose a valuable technique to enhance the distinctiveness of semantics by leveraging exceptionally challenging negative instances, thus enhancing the alignment between visual and linguistic data.

This publication [2] presents a new approach for producing comprehensive descriptions of specific regions within images. The proposed model employs three deep neural networks: one for generating image regions using the Regions with convolutional neural network (R-CNN) technique, another for extracting features using Visual Geometry Group (VGG), and a third for generating descriptive sentences using RNN. The efficiency of the method is validated through experiments conducted on Flickr8K and MSCOCO datasets. The results demonstrate that region-based image descriptions adequately capture the essence of the entire image, often containing even more detailed information than the ground truth sentences. Additionally, the proposed region optimization method enables the selection of suitable regions for an image, achieving descriptive quality comparable to that of full image descriptions. Moreover, the full image description and region-based description can complement each other in conveying information.

In [55] They tackled the issue of learning perception and language together to comprehend the item at a finer level. The use of object descriptions to create a complete knowledge of an object is the central concept of their methodology. They developed two new architectures based on this concept to address related issues such as object captioning and natural language-based object retrieval. While the objective of the object retrieval job is to locate an object given an input query, the objective of the object captioning task is to

simultaneously detect the object and create its related description. They show that hybrid end-to-end CNN-LSTM networks can be used to efficiently handle both issues.

In the study by reference [56], the authors explore the concurrent interdependence of high-level semantic concepts. They propose an innovative approach involving scene-graph-based semantic representation for the task of 'Image Captioning'. This approach stands apart from the existing methodologies for generating captions for images. The image captioning process is deconstructed into two distinct phases: idea cognition and phrase generation. This division facilitates the incorporation of scene graphs as an intermediary stage. The researchers construct a vocabulary comprising semantic terms and advocate for a CNN-RNN-SVM framework to generate the sequence based on scene graphs. This sequential output is then transformed into a binary vector, which serves as input for a subsequent RNN phase. The effectiveness of their technique is assessed using the MS COCO dataset.

In the research conducted by the authors [57], they introduced a language Convolutional Neural Network – Residual Holistic Neural Network (CNN-RHN) model designed specifically for image captioning and well-suited for applications in statistical language modeling. Their innovative language CNN incorporates all preceding words, enabling it to capture significant long-range relationships among historical words. This capability is particularly important for tasks such as 'Image Captioning', distinguishing it from earlier models that relied on a single previous input and hidden state to predict subsequent words. To evaluate the effectiveness of their approach, the researchers employed two datasets, namely Flickr30K and MS COCO.

In this study [5], they suggest a complete deep-learning method for creating 'Image Captions'. Through the use of a semantic attention model, they make use of picture feature information at a specific position every instant to describe the related caption. Through the evaluation of the likeness between sequences of image features and sequences of semantic words, they employed an end-to-end framework for integrating an autonomous recurrent structure as an attention module. Additionally, to enable cross-lingual image captioning, their model is built to transfer the knowledge representation obtained from the English portion into the Chinese portion. They use the most well-known benchmark datasets to

assess the proposed model. On the Flickr8k CN dataset, they show a 3.9% improvement over current state-of-the-art methods for cross-lingual image captioning using the CIDEr metric. The experimental findings show how well their attention Model works.

In the article [58], the authors categorized the current algorithms into two broad groups: top-down and bottom-up approaches. The top-down approach involves directly utilizing the image's information (referred to as visual-level features) to generate a caption, while the latter creates a description using the words that were extracted from the image (known as a semantic-level attribute). To generate image captions, earlier techniques either relied on a one-stage decoder or only used a small portion of visual or semantic-level information. The problem was tackled, and a groundbreaking multi-stage architecture named Stacked visual-semantic (Stack-VS) was introduced. This architecture aimed to generate detailed 'Image Captions' with a focus on both visual and semantic aspects of input images. The approach combined top-down and bottom-up attention models. A novel stack decoder model was proposed, comprising a series of decoder cells. Each cell was equipped with two LSTM layers that collaboratively re-optimized attention weights. These weights were applied to both visual-level feature vectors and semantic-level attribute embeddings, facilitating the creation of intricate image captions. Extensive experiments conducted on the widely recognized MSCOCO dataset demonstrated notable enhancements across various evaluation metrics. In comparison to the leading state-of-the-art method, improvements of 0.372, 1.226, and 0.216 were achieved in BLEU-4, CIDEr, and SPICE scores, respectively.

In this article [59], they suggest a brand-new Pseudo-3D Attention Transfer with a content-aware Strategy (P3DAT-CAS) for the job of captioning images. A pseudo-3D attention network incorporates more detailed information for caption creation and can fully utilize both 2D spatial attention maps and 1D semantic-channel attention features. To produce more elegant phrases, the transfer network is intended to keep attention to situations and provide a wider guide. By choosing captions that are most pertinent to the contents of the images, the content-aware technique can bridge the cross-modal gap between vision and language. The experimental results show that P3DAT-CAS beats

cutting-edge methods. their model outperforms all other models that make use of the cross-entropy training approach in terms of performance.

In this paper [60] the most advanced architectures for sequence modeling tasks like language interpretation and machine translation are transformer-based ones. However, the extent to which they may be applied in multi-modal scenarios like image captioning is still extensively unexplored. they introduce M2 - a Meshed Transformer with Memory for Image Captioning to close this discrepancy. The structure improves both the image encoding and language generation procedures. It integrates learned a priori knowledge to learn a multi-level representation of connections between image regions and uses mesh-like connectivity at the decoding stage to exploit both low- and high-level features. Using an experimental approach, they examined that how the Meshed-memory transformer (M2 Transformer) and several fully attentive models perform compared to recurrent models. Tested on the COCO dataset, their idea demonstrates strong performance in both single-model and ensemble scenarios, specifically on the "Karpathy" test segment and the online evaluation platform. This achievement sets a fresh precedent in the field. They also assess its capability to describe items that were not part of the original training dataset.

In 'Image Captioning' and visual question-answering (VQA), top-down visual techniques have been widely used to facilitate deeper image understanding through fine-grained assessment and even several phases of reasoning. In this paper [61], they suggest a hybrid bottom-up and top-down attention mechanism, allowing for the calculation of attention at the level of objects and other salient picture regions. This is the rationale behind naturally considering attention. In their approach, the allocation of feature weightings is determined by the top-down mechanism, while the bottom-up process (utilizing Faster R-CNN) proposes image regions, each associated with a linked feature vector. When applying our technique to 'Image Captioning,' exceptional results were attained, with scores of 117.9 for CIDEr, 21.5 for SPICE, and 36.9 for BLEU-4 on the MSCOCO test server. This achievement establishes a novel state-of-the-art benchmark. Employing the same methodology to address the VQA problem led us to secure victory in the 2017 VQA Challenge, underscoring the wide-ranging effectiveness of this approach.

They suggest a unique scene-based factoring attention component for image captioning in this study [62]. Their model takes scene concepts into account, in contrast to earlier efforts that either focused on regional features attention or object-centered visual concepts attention. They are the first, as far as they are aware, to consider scene concepts in 'Image Captioning' as well as model relation among scene notions, object-centered visualizations, and caption creation. The LSTM hidden feature in this proposed scene-based factored attention module directly incorporates scene notions in the form of a factored tensor. they determined the relative importance of regional features and object-centered visual concepts depending on the concealed feature that is incorporated into the scene. The true strength of our suggested module is its capacity to pay attention to hierarchically visual data for improved captions. The effectiveness of the suggested approach has been validated by experiments using the MS COCO captioning datasets.

In this study [63] A method for generating medical reports for retinal images is presented, utilizing expert-defined keywords and an innovative attention-based strategy. This method can predict technical keywords and combine them for advanced word sampling. Experimental results demonstrate that the proposed model produces more accurate and meaningful descriptions for retinal images, with a performance increase of approximately 74% in BLEU average, 63% in ROUGE, 87% in CIDEr, and 63% in METEOR compared to non-keyword attention-based baselines. Attention visualization reveals intriguing patterns of potential symptoms in specific image regions. To further enhance the explainability of ML-based models for retinal image captioning, suggesting an automatic metric for measuring explainability is a promising avenue for future research within our community.

In this study [64], they introduce a straightforward yet powerful prompt-driven framework for generating image captions, a concept that has received limited attention in the captioning community. Through prompt engineering, our approach achieves the generation of captions exhibiting diverse styles. To delve deeper into the capabilities of prompt learning, they guide the network to autonomously discover appropriate prompt vectors within the continuous word embedding space. Their qualitative and quantitative experiments substantiate the efficacy of the proposed framework. They attain results on

two distinct image captioning benchmarks, namely the COCO Karpathy split and TextCaps, employing a consolidated model.

In this study [65], they explore the enhancement of visual-language alignment and linguistic coherence in a diffusion model designed for image captioning. To substantiate our assertion, they introduce a novel Semantic-Conditional Diffusion Networks (SCD-Net) semantic-conditional diffusion process that enriches the existing diffusion model with additional semantic priors. Additionally, they implement a guided self-critical sequence training strategy to stabilize and enhance the diffusion process. Their proposed approach is empirically demonstrated to outperform state-of-the-art non-autoregressive methods. Notably, their new diffusion model paradigm exhibits superior performance compared to a competitive autoregressive method with the same Transformer encoder-decoder structure. The findings underscore the promising potential of the diffusion model in the context of image captioning.

## 3.3 Literature Review Results Analysis

All researchers complied that a relevant problem in this research domain is the issue of accuracy and performance in the papers. The major gap found after the literature review is the lack of accuracy in the generated captions and the gap is due to the lack of semantic knowledge. All research tries to solve this mentioned problem by performing different experiments with different machine learning and deep learning techniques to improve the performance of the image caption. After the experiment, all show the quantitative result of different evaluation metrics to evaluate the accuracy and performance of generated image caption. Table 3.1 shows all the research techniques and methods that are used for image captioning and the results of different evaluation metrics are also shown in Table 3.1.

**Table 3.1** Analysis of Findings from the Literature Review

| Sr# | Ref | Year | Author | Model | B1 | B2 | B3 | B4 | M | C | R | S |
|-----|-----|------|--------|-------|-----|-----|-----|-----|-----|-----|-----|-----|
| 1 | [63] | 2023 | Ting-Wei Wu | ML-based medical report generation system | 69.69 | 61.95 | 54.96 | 50.08 | 70.44 | 56.50 | 72.52 | - |
| 2 | [64] | 2023 | Ning Wang | ConCap | - | - | - | 40.5 | 30.9 | 133.7 | - | 23.8 |
| 3 | [65] | 2023 | Jianjie Luo | SCD-Net | 80.2 | 64.9 | 50.1 | 38.1 | 29.0 | 126.2 | 58.5 | - |
| 4 | [49] | 2022 | Yang Yang | I-CPRA | - | - | - | - | 27.9 | 111.2 | - | 21.0 |
| 5 | [19] | 2022 | Muhammad Abdelhadie Al-Malla1 | YOLO | 6.26 % | 8.42 % | 11.53 % | 16.09 % | 3.82 % | 15.04 % | 3.76 % | 5.88 % |
| 6 | [20] | 2021 | HAIYANG WEI | Visual Attention model | 80.5 | 65.7 | 51.0 | 38.9 | 28.3 | 126.7 | 58.8 | 21.7 |
| 7 | [50] | 2021 | Junlong Feng | SF-GAN | 27.9 | 14.7 | 8.3 | 4.8 | 11.2 | 42.2 | 26.2 | 15.0 |
| 8 | [53] | 2021 | Kien Nguyen | SG2Caps – RL | | | | 33.0 | 26.2 | 112.3 | 55.3 | 19.4 |
| 9 | [46] | 2020 | Beichen Zhang | SSAAL | - | - | - | 34.3 | 26.2 | 106.2 | 55.3 | - |

| 10 | [47] | 2020 | Fenglin Liu | Textual Distilling Module and Textual Association Module | 80.9 | 65.7 | 51.2 | 39.3 | 29.5 | 129.0 | 59.2 | - |
|----|------|------|-------------|------------------------------------------------------------|------|------|------|------|------|-------|------|---|
| 11 | [48] | 2020 | Seung-Ho | a domain-specific image caption generator | 79.1 | 62.4 | 47.1 | 35.9 | 27.1 | - | - | - |
| 12 | [54] | 2020 | Wenjie Cai | image captioning framework | 80.8 | 64.3 | 49.6 | 37.5 | 28.2 | 126.0 | 58.2 | 21.8 |
| 13 | [5] | 2020 | Bin Wang | end-to-end deep learning approach | 66.8 | 46.8 | 32.2 | 22.1 | - | 55.12 | 20.36 | - |
| 14 | [58] | 2020 | Ling Cheng | Stack-VS | 79.0 | 63.4 | 48.9 | 37.2 | 28.8 | 118.9 | 57.5 | - |
| 15 | [60] | 2020 | Cornia, Marcella | M2 Transformer | 81.6 | 66.4 | 51.8 | 39.7 | 29.4 | 129.3 | 59.2 | - |
| 16 | [45] | 2019 | Chen He | VSDA | 75.3 | 59.1 | 45.1 | 34.4 | 26.5 | 53.2 | 55.2 | - |

| 17 | [51] | 2019 | Zongjia n Zhang | (FCN)-LSTM | 71.2 | 51.4 | 36.8 | 26.5 | 24.7 | 88.2 | - | - |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 18 | [3] | 2019 | Kunpen g Li | VSRN | 71.3 R@1 | 90.6 R@5 | 96.0 R@10 | - | - | - | - | - |
| 19 | [55] | 2019 | Anh Nguyen | hybrid end-to-end CNN-LSTM networks | 69.6 | 57.0 | 46.9 | 40.0 | 34.2 | 163.2 | 68.5 | - |
| 20 | [59] | 2019 | Jie Wu | P3DAT-CAS | 75.3 | 59.0 | 45.0 | 34.1 | 27.0 | 109.1 | 55.4 | - |
| 21 | [62] | 2019 | Chen Shen | scene-based factored attention module | 80.3 | 64.6 | 60.1 | 38.1 | 28.5 | 126.8 | 58.2 | 22.0 |
| 22 | [6] | 2018 | Lizhao Gao | Scene graph (CNN-RNN-SVM) framewor k | 67.6 | 49.3 | 35.5 | 26.1 | 22.3 | 76 | - | - |
| 23 | [61] | 2018 | Peter Anderso n | Up-down | 80.2 | 64.1 | 49.1 | 36.9 | 27.6 | 117.9 | 57.1 | 21.5 |
| 24 | [57] | 2017 | Jiuxiang Gu | CNN+RH N | 72.3 | 55.3 | 41.3 | 30.6 | 25.2 | 98.9 | - | - |
| 25 | [52] | 2016 | Quanzen g You1 | ATT-FCN | 70.9 | 53.7 | 40.2 | 30.4 | 24.3 | - | - | - |

| 26 | [2] | 2015 | Xiaodan Zhang | R-CNN, VGG and RNN | 64.00 | 44.57 | 24.87 | 12.54 | - | - | - | - |
|----|-----|------|---------------|--------------------|-------|-------|-------|-------|---|---|---|---|

## 3.4  Summary

In this chapter, the literature review of this research domain 'Image Captioning' was discussed. In the literature review, the problem and gap in the research and different techniques and methods used to solve the problem of this research domain were discussed. The results of different authors that come after experimenting using different machine learning and deep learning techniques are also discussed.
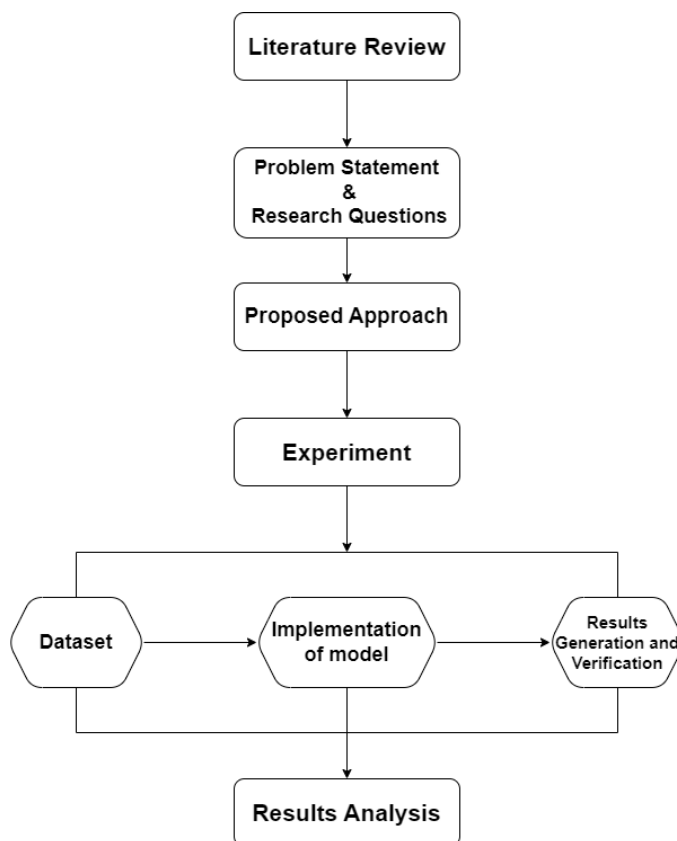
# Chapter 4

# Research Methodology and Proposed Framework

The main objective of this chapter is to describe the research methods or procedures used in the conduction of this research. Image captioning is the main area of study in artificial intelligence. Although various studies on image captioning have been carried out in the past, it is still difficult to produce appropriate visual descriptions for an image. All the measures implemented that outlined in Figure 4.1 to improve the effectiveness of 'Image Captioning'. My proposed approach to conducting this research will also be discussed. The experiment stage of this research will also be discussed and how practical performed during this research and evaluation metrics to evaluate the results.

## 4.1   Research Methodology for Image Captioning

A methodology is referred to as a process or a set of rules or concepts. It might offer guidelines from which precise methods or procedures could be easily understood, used to solve problems within the domain, and changed to address issues within the domain's defined boundaries. The research problem was identified clearly in this research and after this proceeded this research to solve that research problem. there are following steps for the research completion all steps are discussed below:

**Figure 4.1** Research Methodology for Image Captioning

**Literature Review:** The research area/domain in which research conducted me, 'Image Captioning,' was selected first. This field is an important research area of computer vision and natural language processing. After selecting the research domain, multiple research papers were studied to find out the research gap and problem. All papers related to 'Image Captioning' in which research used different research techniques or methods to perform research in 'Image Captioning'.

 **Problem Identification:** In the second step, after studying the research papers. The lack of accuracy and performance in the 'Image Captioning' field was found to be a research gap. Accuracy is reduced due to the lack of semantic attention in the processing of images because of relationships or connections among the objects within the image.

**Experiment:** In the next step, after problem identification, my research proceeds towards the experiment to conduct my research to solve the mentioned research problem. An encoder-decoder approach is proposed to fill the research gap, after performing the

experiment result is calculated and verified by using the evaluation metrics used for the evaluation of generated image captions.

**Results Analysis:** In the last step, after completing the experiment both quantitative and qualitative results and analysis were performed. Compare all the results that come from the evaluation metrics with the results of the other research perform in the field of 'Image Captioning'. Captions generated from my proposed approach compared with captions generated from the other research approaches.

**Base Paper:** To compare my research work perform by H. Wei [19], [20] research article is followed as a base paper. They used a visual attention model for a generation of image captioning and focused on the semantics of the image while creating a caption for an image [19], [20]. They used the same evaluation metrics like BLEU-1, BLEU-2, BLEU-3, BLEU-4, CIDEr, METEOR, and ROGUE-L, the results of their experiment are 80.5, 65.7, 51.0, 38.9, 28.3, 126.7 and 58.9 respectively. Results compared with the results of base paper results.

## 4.2   A Proposed Approach: ViT-GPT-2

To overcome the gap of image caption found through the literature review. An encoder-decoder framework is proposed to solve the research gap. The major issue in the 'Image captioning' is the accuracy in the generated image captions. Accuracy is low due to a lack of semantic knowledge in 'Image Captioning'. This encoder-decoder framework mainly focuses on the semantic knowledge to improve the accuracy of 'Image Captioning'. Semantic knowledge gained through the enhancement of the feature extraction. All the objects and the connections among the objects are extracted accurately. A transformer-based technology was proposed to improve the accuracy of the image caption. The proposed approach is shown in Figure 4.2.

**Encoder:** In the proposed transformer-based encoder-decoder approach the Vision transformer (ViT) works as an encoder. The detailed working of the Vision Transformer is discussed in Chapter 2 Section 2.10. The Vision Transformer (ViT) is a ground-breaking neural network design that extends the Transformer to the field of computer vision. The

Vision Transformer is used to extract features for the image, the image is divided into patches and each patch moves forward to the vision transformer for further processing. The number of patches increased by changing its pixel ratio from 16x16 to 8x8. For creating patches of an image an example is discussed below:
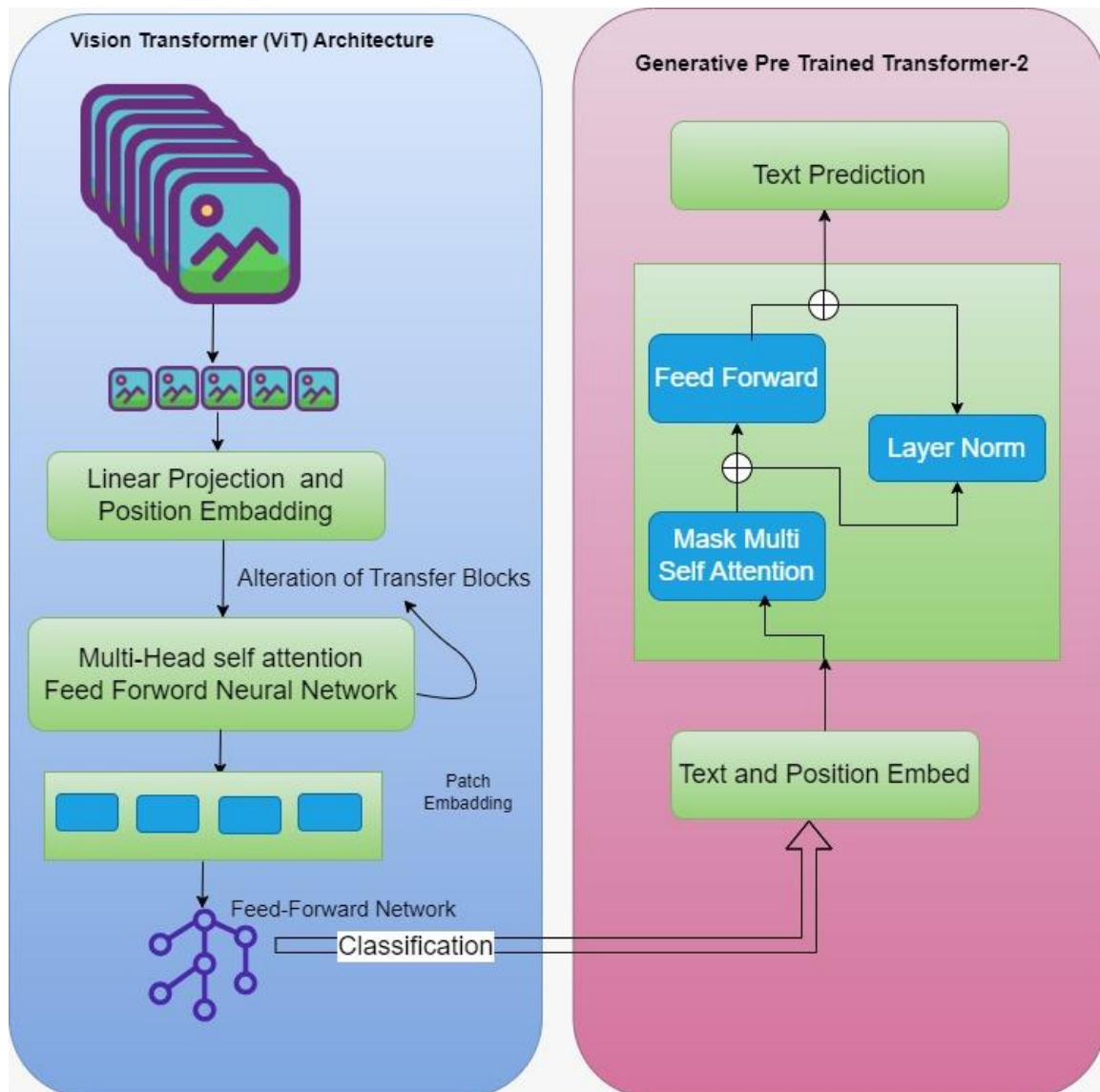
Divide the image of 240 X 240 pixels to create the patches of 16 X 16 pixels and 8X 8 pixels to find the number of patches given in equations 4.1 and 4.2 respectively.

$$Number\ of\ patches = \frac{image\ width}{patch\ width}\ X\ \frac{image\ height}{patch\ height} \tag{4.1}$$

$$Number\ of\ pactches\ with\ 16\ X\ 16 = \frac{240}{16}\ X\ \frac{240}{16} = \frac{57600}{256} = 225 \tag{4.2}$$

$$Number\ of\ pactches\ with\ 8\ X\ 8 = \frac{240}{8}\ X\ \frac{240}{8} = \frac{57600}{64} = 900 \tag{4.3}$$

As in this research major focuses on the exploitation of semantic knowledge, the number of patches increases so that feature extraction and semantic knowledge from each patch of the image. ViT conceptualizes images as sequences of patches, each of which goes through an embedding procedure to generate a token, rather than processing images pixel by pixel. By breaking down images into patches and converting them into embeddings, ViT effectively captures both local and global features. These embeddings encode spatial relationships and visual context, providing a comprehensive representation of the image. Leveraging self-attention mechanisms, ViT captures intricate dependencies between patches, enhancing its ability to discern important visual cues.  These tokenized patches are used as the input to the transformer encoder, which is enhanced with positional embeddings to preserve spatial information working of the Vision transformer as shown in Figure 4.2.

**Figure 4.2** ViT-GPT-2 Encoder-Decoder Framework

**Decoder:** In the proposed transformer-based encoder-decoder approach the Generative Pre trained Transformer 2 (GPT-2) works as a decoder in this framework. The GPT-2 decoder role in image captioning occurs after an image encoder processes the visual content extracted from the ViT encoder. GPT-2 produces comprehensible and contextually pertinent subtitles that go along with the visual content by training the decoder on both visual and textual information. GPT-2 generates captions word by word based on its understanding of the previously generated text and the visual content of the image. By integrating visual and semantic attention mechanisms, GPT-2 focuses its attention on both

the evolving textual context and salient image features, ensuring that the generated words harmonize with the image's content and the evolving caption. The quality of the generated captions is improved by fine-tuning GPT-2 for image captioning, which improves its capacity to align and ground textual output to the visual elements.

## 4.3   Experiment

In this section, the complete process to experiment, as shown in Figure 4.3, will be discussed. First of all, take a dataset of images and preprocess that dataset, after preprocessing dataset data move forward towards encoder (ViT). Feature extracted using an encoder and its output as input in decoder (GPT-2). In the last step of this experiment, the captions generated are evaluated using evaluation metrics that are used for the image captions. After these results are compared with the results of other researchers who perform in the field of 'Image Captioning'

**Figure 4.3** Experimental Study for Image Caption

## 4.3.1   Preparing the data

Several datasets are used to test, train, and evaluate the 'Image Captioning' systems. The datasets differ from one another in a variety of areas, such as the number of images,

and the number of descriptions for each picture, and these captions are organized according to the images. There are three widely used datasets: Flickr8k [66], Flickr30k [67], and MS COCO Dataset [68].

To perform the complete experiment, first of all, A dataset with a large number of images and captions relevant to those images is selected by me. Each image has a one-to-one mapping with the textual caption regarding those images. There are some dataset exits that match the criteria like MS COCO and any other data set that fulfill all the requirements of the dataset like a large number of images and captions that are related to that image. To perform this experiment, the MS COCO dataset was used for the image captioning process. A detailed description of the dataset used for the experiment is given below:

**COCO Dataset:** The primary objective of the COCO Object Detection Task is to push the boundaries of object detection advancements. COCO presents two distinct tasks for object detection, with one utilizing bounding box outputs and the other involving object segmentation outputs, which is occasionally referred to as instance segmentation. It is possible to utilize the COCO train, validation, and test collections, encompassing over 200,000 images across 80 distinct object categories. Each object instance is accompanied by a comprehensive segmentation mask. The training and validation sets include publicly available annotations, consisting of over 500,000 segmented object instances [68].

**Flicker30k Dataset:** The Flickr30K Dataset, known as a resource for facilitating automated image description and enhancing language comprehension, is comprised of 30,000 photographs sourced from Flickr, each accompanied by 158,000 human-generated captions [67]. Notably, this dataset does not come with predetermined partitions for training, testing, or validation purposes. Instead, researchers have the flexibility to customize their own sets for validation, testing, and training as needed. Furthermore, the dataset incorporates classifiers for identifying colors, detectors for recognizing common objects, and a preference for selecting larger objects within the images.

### 4.3.2  Pre-processing

The combination of ViT and GPT-2 in an image captioning task involves two separate pre-processing pipelines: one for the image data preprocessed ViT and one for the text data preprocessed by GPT-2.

**Resizing:** After gathering images in the form of a dataset need to fix the size of the image. In this research for ViT 240x240 pixels size images to ensure uniformity in the process.

**Patch Extraction:** After resizing the image to a fixed ratio. Each image is converted into patches of 8x8 pixel from 240x240 pixel image and every patch is fixed ratio square and non-overlapping patch. Each separate patch is treated as a separate input in ViT.

**Tokenization:** After gathering dataset images and their relevant image captions, all captions are tokenized into individual words or sub words, ViT-GPT-2 processes text into the tokenized form.

**Special Tokens:** Add special tokens to tokenized captions such as the start of the sequence and end of the sequence that helps the model to understand the start and end of the caption. A padding token is also used to ensure the length of the caption.

**Noise Removal:** Our dataset comprises two components: an image dataset and a set of captions corresponding to those images. Within the captions list, some entries are either missing or inconsistent. To mitigate this noise, a systematic process is employed to review each caption individually. Any irrelevant or empty captions identified during this review are subsequently removed from the list.

### 4.3.3  Feature Extraction

In image captioning tasks that involve a combination of ViT and GPT-2, the visual features are typically extracted by the ViT model, while GPT-2 is used for processing and

generating textual descriptions. The ViT model extracts features from the image, and these visual features are then combined with the textual features from GPT-2 for further processing. Features extracted by the vision transformer is equal to the number of patches given to the transformer [39] and all features uses as an input by the decoder GPT-2.

ViT extracts all visual features and spatial information from each of the patches of the image. GPT-2 processes textual information such as image captions and captures semantic information from the image captions

## 4.3.4  Encoder

Image captioning is one of many computer vision tasks for which the Vision Transformer (ViT) works as an encoder. The ViT processes input images as an encoder and converts them into useful representations used by ViT as follows:

**Image Patching:** The vision transformer takes the preprocessed image from the image dataset and divides the image into small patches according to their pixels, images divided into 8x8 pixel patches. All patches are non-overlapping patches and every patch is treated as a separate input image to the vision transformer.
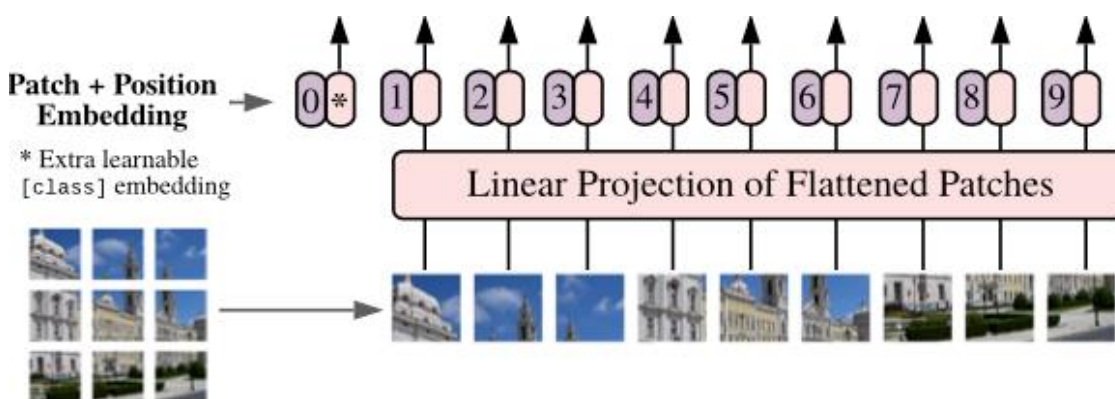


**Figure 4.4** Presentation of image patching

**Patch Embedding:** After making the patches of an image each patch is converted into embedding using the linear projection layer of the vision transformer. All of these patches have both semantic and visual information. To feed the patches into the ViT model, the

patches are reshaped into a linear format. By doing this, pixel data is transformed into a feature vector.

**Positional Encodings:** After embedding the patches in the vision transformer positional encoding is added to the patch embeddings to find out the semantic and spatial information from the patch of the image. This process helps the model to understand the arrangements of the patches of an image. Trigonometric functions are frequently used to create positional encodings to capture positional connections.

The initial representations for each patch are the patch embeddings.



**Figure 4.5** Presentation of patch Embedding

**Transformer Encoder Layers:** Each patch embedding serves as the input to the Transformer encoder layers (after positional encodings are added). These embeddings show the visual characteristics of the image at various spatial scales.

A. **Multi-Head Self Attention:** After positional encoding both positional encodings and patch embeddings are input to various layers of self-attention. The model can recognize relationships between various patches as well as recognize their context within the full image by using self-attention. Based on the importance between positions, self-attention calculates weighted sums of values for each position in the input sequence. Each attention head gains the ability to concentrate on different components of the input sequence, enabling the model to recognize various kinds of relationships.

B. **Feed-Forward Nural Network**: The patch embeddings go through position-wise feed-forward neural networks after self-attention. Fully interconnected layers with activation capabilities (like ReLU) make up this network. By taking into account interactions between patches, these networks can learn more features and properties of the image through the image patches.

C. **Normalization and Residual Connections:** After the self-attention and feed-forward network operations for each point independently, layer normalization is used. The vanishing gradient issue is avoided by using residual connections to include the original input embeddings in the normalized outputs.

**CLS Token Representation:** The whole semantic knowledge of the image is captured by the representation of the [CLS] token, which is frequently prepended to the input sequence. This representation serves as the decoder's input for the retrieved visual feature.

The generation of a condensed portrayal capturing the entirety of the image's content is achieved by the ViT encoder after processing the patches and encoding spatial and contextual relationships. Then, a variety of tasks, including object detection, captioning, and classification, can be carried out using this representation.

## 4.3.5 Decoder

The image captioning process uses GPT-2 as the decoder to provide evocative captions based on the visual features that are extracted from the encoder vision transformer. GPT-2 works as a decoder as follows:

**Image Encoding:** An image encoder (such as Vision Transformer) first processes the input image to create a fixed-size feature vector that reflects the visual content of the image.

**Caption Starts Token:** The initial hidden state for the GPT-2 decoder is the [CLS] token representation from the ViT's output. This token gives the decoder a context in which to produce captions.

**Language Modelling:** By predicting the following token in the sequence based on the previous ones, the GPT-2 decoder creates captions. The initial state of the decoder is

determined by the [CLS] token representation, and as decoding proceeds, the model pays attention to pertinent areas of the visual feature to produce contextually relevant text.

**Self-Attention and Feed Forward Layers:** GPT-2 employs self-attention mechanisms at each stage of decoding to take into consideration the connections between the generated tokens and the visual feature representation. This aids the model in producing captions that are cohesive and meaningful in terms of semantics and the content of the image. These layers learn to produce content that is coherent and contextually appropriate by capturing the relationships between words in context.

**Text Generation:** The GPT-2 decoder creates tokens one at a time, predicting the following tokens from the previous ones using its learned language model. Until an end-of-sequence token or a maximum sequence length is reached, this operation is repeated.

The image captioning system makes use of both models' strengths by combining ViT as the encoder to extract visual information and GPT-2 as the decoder to produce text. GPT-2 creates descriptive captions based on the visual elements that ViT extracts from the image. By using this strategy, the system provides comprehensive captions, pertinent to the context, and accurately reflects the content of the input image.

## 4.4   Evaluations Metrics

Within the realms of computer vision and natural language processing, generating descriptions for images, known as 'Image Captioning', poses a considerable challenge. Reliable measures that can gauge how closely generated captions match reference captions created by humans are needed to assess the quality of generated captions. Bleu, CIDEr, Rogue, and meteor are image captioning evaluation metrics that are often employed.

These approaches compute the scores of the predictions using certain widely used metrics, including Precision, Recall, and F1 Score. The True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN) metrics are used to assess performance. As indicated below, the TP, TN, FP, and FN values are discovered by creating a confusion matrix and comparing the actual and anticipated values.

The values for Precision, Recall, and F1 Scores using the TP, TN, FP, and FN data were determined.

**Precision:** The percentage of successfully predicted positive values among all accurately predicted positive values.

$$Percision = \frac{TP}{TP+FP} \tag{4.4}$$

**Recall:** Determines the proportion of correctly predicted positive values among all positive values.

$$Recall = \frac{TP}{TP+FN} \tag{4.5}$$

**F1-Score:** Precision and Recall are weighted average.

$$F1 - Score = 2 * \frac{Recall*Percision}{Recall+Percision} \tag{4.6}$$

**BLEU (Bilingual Evaluation Understudy):** BLEU is a well-liked statistic for comparing automatically generated captions against reference captions. Between the created caption and the reference captions, it calculates the n-gram overlap. Better alignment between the prepared and source captions is indicated by higher BLEU ratings [69]. The effectiveness of the BLEU metric is influenced by two factors: the quantity of reference translations and the length of the text. N-grams are used in the calculation of the updated precision metric that [69] presented. The reason BLEU is well-known is that it was a pioneer in the field of automatic evaluation of machine-translated text and that its results have a good connection with human evaluations of quality [70], [71].

Let's look at an example to better understand how the BLEU 1-4 scores are determined. Following is the formula for calculating precision for n-gram.

$$Precision\ n - gram = \frac{Number\ of\ correct\ predicted\ n-gram}{Number\ of\ total\ predicted\ n-grams} \tag{4.7}$$

The calculation of Geometric Average Precision Scores for individual n-grams involves utilizing the formula presented in equation 4.5. This computation follows the determination of Precision for the respective n-grams. Equation 4.6's Brevity penalty, when

added to the Geometric Average Precision as previously indicated, yields the BLEU-N score, which is shown in equation 4.7.

$$Geometric\ Average\ Percison\ Score\ (N) = exp(\sum_{N=1}^{N} W_n logp_n) = \prod_{n=1}^{N} P_n^{W_n}$$

(4.8)

$$Brevity\ Penalty = f(x) = \begin{cases} 1, & if\ c > r \\ e^{\left(1-\frac{r}{c}\right)}, & if\ c \leq r \end{cases}$$

(4.9)

The variable 'r' represents the length of the reference sentence, while 'c' indicates the length of the predicted sentence.

$$BLEU(N) = Brevity\ Penalty.\ Geometric\ Average\ Precision\ Score(N)$$

(4.10)

**METEOR (Metric for Evaluation of Translation with Explicit ORdering):** METEOR stands as a distinct metric, evaluating the quality of produced captions by considering both phrase- and unigram-based similarity, along with a penalty component for variations in word order. It tries to reward accuracy and fluency in the captions that are created [72]. Referential texts and common word sections are contrasted. Additionally, sentence stems and word substitutes are taken into account while matching. METEOR can produce superior segment or sentence-level correlation.

$$METEOR\ SCORE = \frac{10*Percision*Recall}{Recall+9Percision}$$

(4.11)

**CIDEr (Consensus-based Image Description Evaluation):** The CIDEr measure considers the agreement between the reference captions for an image. It takes into account both n-gram matches and the significance of those matches in the source captions. CIDEr assists in rewarding distinctive and varied captions that are both fluent and imaginative [73]. It is a machine-learning-based consensus metric for assessing picture descriptions.

Frequently, existing datasets typically contain merely five captions for each image. Prior evaluation instruments can gauge agreement between machine-generated captions and human evaluations, but this is restricted by the small dataset size. Nevertheless, CIDEr utilizes the concept of term frequency-inverse document frequency to acquire human consensus, as explained in [74].

$$g_k(s_{ij}) = \frac{h_k(s_{ij})}{\sum_{\omega_l \in \Omega} h_l(s_{ij})} \log\left(\frac{|I|}{\sum_{I_p \in I} \min\left(1, \sum_q h_k(s_{pq})\right)}\right) \qquad (4.12)$$

$$\text{CIDEr}_n(c_i, S_i) = \frac{1}{m} \sum_j \frac{g^n(c_i) \cdot g^n(s_{ij})}{\|g^n(c_i)\| \|g^n(s_{ij})\|} \qquad (4.13)$$

$$\text{CIDEr}(c_i, S_i) = \sum_{n=1}^{N} w_n \text{CIDEr}_n(c_i, S_i) \qquad (4.14)$$

**SPICE (Semantic Propositional Image Caption Evaluation):** A statistic called SPICE assesses how semantically comparable the created and reference captions are. Semantic parsing is used, and the effectiveness of the resulting captions is assessed using semantic statements [75]. It depends on a semantic model of graphs known as a scene graph [76], [77]. This graph may take image captions and extrapolate information about various items, properties, and their relationships.

**ROUGE (Recall-Oriented Understudy for Gisting Evaluation):** ROUGE is a statistic that was initially created to assess text summarization. By taking into account the overlap between created and reference descriptions based on n-gram matching, it may also be modified to assess image captioning. ROUGE exists in several iterations, encompassing ROUGE-1, ROUGE-2, ROUGEW, and ROUGE-SU4. Specifically, ROUGE-1 and ROUGE-W are applicable when evaluating individual documents, whereas ROUGE-2 and ROUGE-SU4 exhibit effectiveness in the context of summarizations. Nevertheless, the evaluation of multi-document text summaries poses challenges for the ROUGE metric [78].

$$\text{ROGUE} - \text{L F1} = 2 * \frac{Recall*Percision}{Recall + Percision} \qquad (4.15)$$

## 4.5  Summary

In this chapter research methodology about my research work was discussed and discuss my proposed work to solve the research problem and to achieve research objectives. The complete experimental setup of my proposed framework also discussed in this chapter.

# CHAPTER 5

# RESULTS

In this chapter, the result of this research that comes from the experiment was discussed. In the final step, the performance of my research approach with the other approaches that have already been used for the research in the same domain is also discussed.

My research is conducted in the domain of 'Image Captioning'. The major research gap is the lack of accuracy in the caption generated for an image. This lack of accuracy is due to the lack of semantic knowledge in the area of 'Image Captioning'. To solve this mentioned problem, An experiment was performed using an encoder-decoder transformer-based approach by using ViT and GPT-2. MS COCO and Flicker30k datasets were used to experiment. After experimenting, the results were evaluated on the specific evaluation metrics that are used for 'Image Captioning'.

## 5.1 Results of Image Captioning with ViT-GPT-2

In this research, the major problem is the lack of accuracy in the generated captions and this lack of accuracy is due to the lack of semantic knowledge of an image. Our main focus is to improve the performance of the 'Image Caption' that is generated from the

different objects and features of an image and explore the semantic knowledge of the image and the relationships of the objects within the image.

To achieve research objectives an encoder-decoder-based framework was proposed. In this framework, the vision transformer (ViT) is used as the encoder that extracts the features from the and sends output to GPT-2 used as the decoder that generates the captions for images in this research. Further in this section, the result of our work will be discussed and will be compared with the research work in this domain.
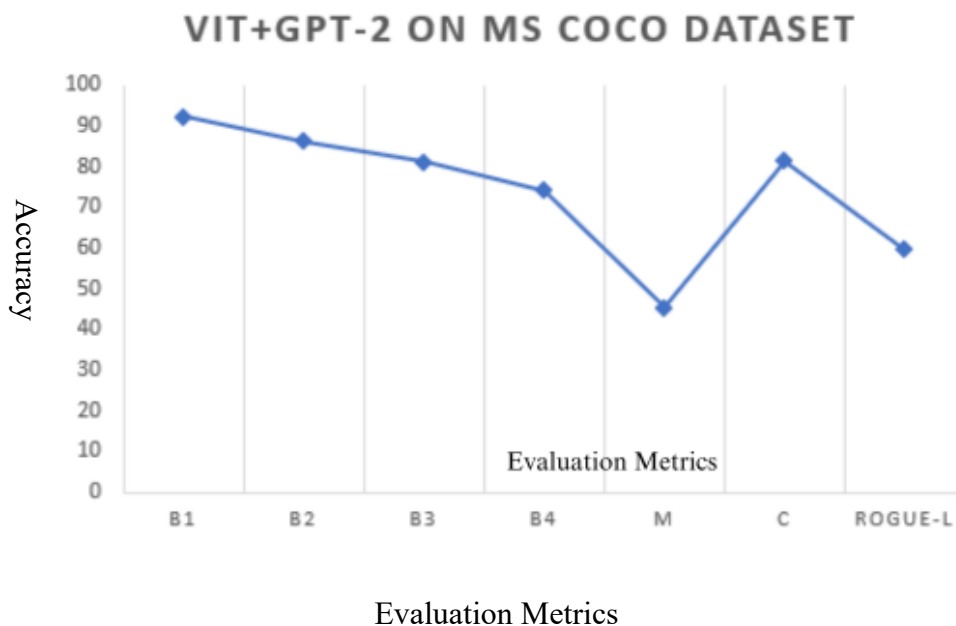
### 5.1.1  Quantitative Analysis

Table 5.1 shows the results that are generated from our model ViT-GPT-2. These results are compiled from the MS COCO Dataset. Similarly, Table 5.2 shows the results generated from our model ViT-GPT-2 on the flicker 30k dataset.

Reliable measures that can gauge how closely generated captions match reference captions created by humans are needed to assess the quality of generated captions. Bleu-1 (B1), Bleu-2 (B2), Bleu-3 (B3), Bleu-4 (B4), CIDEr C, Rogue, and Meteor M are image captioning evaluation metrics are often employed to evaluate the generated image captions.

**Table 5.1** Result of ViT-GPT-2 on MS COCO Dataset

| Model | B1 | B2 | B3 | B4 | M | C | Rogue-L |
|---|---|---|---|---|---|---|---|
| ViT-GPT-2 | 92.18 | 86.15 | 81.17 | 74.19 | 45.36 | 81.37 | 59.5 |

**Figure 5.1** Graphical Representation of Results on MS COCO Dataset

**Qualitative Result Analysis on MS COCO Dataset:**

**BLEU Scores:** BLEU measures the overlap between the generated captions and reference captions in terms of n-grams such as bleu-1, bleu-2, bleu-3, and bleu-4. A higher bleu score indicates the similarity in the generated image caption with the reference caption. In this research, all of the bleu scores improved as compared to the other research. Results from the experiment show a performance improvement as Bleu-1, Bleu-2, Bleu-3, and Blue-4 at 92.18, 86.15, 81.17, and 74.19 respectively.

**METEOR:** METEOR considers precision, recall, and alignment of n-grams, providing a more holistic view of performance. Meteor score shows the quality of the generated caption in the context of the image caption. In this research meteor score is 45.36, this score is improved from other research shows that the quality of the generated image caption improves.
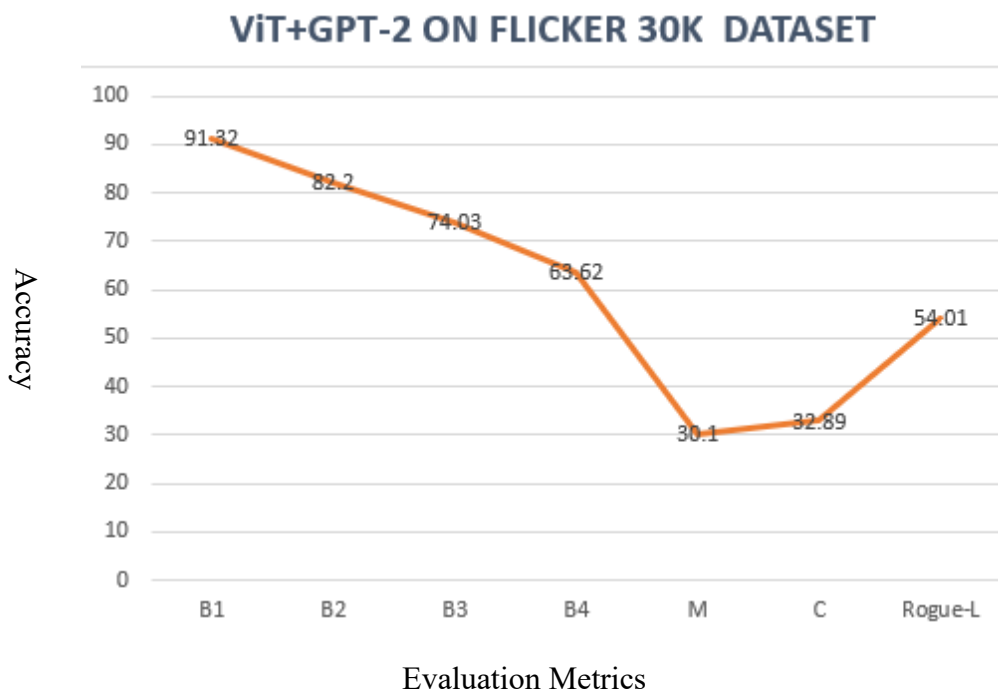
**CIDEr:** CIDEr emphasizes the importance of generating diverse and descriptive captions. Higher CIDEr scores indicate better performance in capturing the essence of the images. In this research, our research model gained a CIDEr score is 81.37.

**ROUGE-L:** ROUGE-L measures the overlap of Longest Common Subsequences (LCS) between generated and reference captions. A higher ROUGE-L score means better recall of important content in the captions. The score of ROUGE-L from this research is 59.5 this also improves from the other research shows that all important information is extracted from the image using the ViT-GPT-2 framework.

As

**Table 5.2** Result of ViT-GPT-2 on Flicker 30k Dataset

| Model | B1 | B2 | B3 | B4 | M | C | Rogue-L |
|-------|------|-------|-------|-------|------|-------|---------|
| ViT-GPT-2 | 91.32 | 82.20 | 74.03 | 63.62 | 30.1 | 32.89 | 54.01 |



**Figure 5.2** Results Graphical Representation of Flicker 30k Dataset

**Qualitative Result Analysis on Flicker 30k Dataset:**

**BLEU Scores:** In this research, all of the BLEU score also improves with the flicker 30k dataset as compared to the other research. Results from the experiment show a

performance improvement as Bleu-1, Bleu-2, Bleu-3, and Blue-4 at 91.32, 82.20, 74.03, and 63.62 respectively.

**METEOR:** In this research meteor score on the Flicker 30k dataset is 30.1, this score is improved from other research shows that the linguistic quality of the generated image caption improves.

**CIDEr:** Higher CIDEr scores indicate better performance in capturing the essence of the images. In this research, our research model gained a CIDEr score is 32.89 on the Flicker 30k dataset.

**ROUGE-L:** The score of ROUGE-L from this research is 54.01 on the Flicker 30k dataset, this also improves from the other research shows that all important information extracted from the image using ViT-GPT-2 framework.

## 5.1.2 Comparison Of Our Results With Other Models On MS COCO Dataset

Table 5.3 shows the result of the experiment on the MS COCO dataset with the comparison of other models with our model that major focus on the semantic knowledge of the image during the generation of the image captioning. Different model works differently and different results are generated from the experiment.

**Table 5.3** Result Comparison with Other Models on MS COCO Dataset

| Model Name | B1 | B2 | B3 | B4 | M | C | Rogue-L |
|---|---|---|---|---|---|---|---|
| **ML-based medical report System** [63] | 69.69 | 61.95 | 54.96 | 50.08 | 70.44 | 56.50 | 72.52 |
| **ConCap** [64] | - | - | - | 40.5 | 30.9 | 133.7 | - |
| **SCD-NET** [65] | 80.2 | 64.9 | 50.1 | 38.1 | 29.0 | 126.2 | 58.5 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Visual Attention Model** [20] | 80.5 | 65.7 | 51.0 | 38.9 | 28.3 | 126.7 | 58.8 |
| **VSDA** [45] | 75.3 | 59.1 | 45.1 | 34.4 | 26.5 | 53.2 | 55.2 |
| **SSAAL** [46] | - | - | - | 34.3 | 26.2 | 106.2 | 55.3 |
| **Textual Distilling Module and Textual Association Module** [47] | 80.9 | 65.7 | 51.2 | 39.3 | 29.5 | 129.0 | 59.2 |
| **a domain-specific image caption generator** [48] | 79.1 | 62.4 | 47.1 | 35.9 | 27.1 | - | - |
| **I-CPRA** [49] | - | - | - | - | 27.9 | 111.2 | - |
| **SF-GAN**$_{RF}$ [50] | 27.9 | 14.7 | 8.3 | 4.8 | 11.2 | 42.2 | 26.2 |
| **(FCN)-LSTM** [51] | 71.2 | 51.4 | 36.8 | 26.5 | 24.7 | 88.2 | - |
| **YOLO** [19] | 6.26 % | 8.42 % | 11.53 % | 16.09 % | 3.82 % | 15.04 % | 3.76 % |
| **ATT-FCN** [52] | 70.9 | 53.7 | 40.2 | 30.4 | 24.3 | - | - |
| **VSRN** [3] | 76.2@R1 | 94.8R@5 | 98.2@R10 | - | - | - | - |
| **SG2Caps** [53] | | | | 33.0 | 26.2 | 112.3 | 55.3 |
| **image captioning** | 80.8 | 64.3 | 49.6 | 37.5 | 28.2 | 126.0 | 58.2 |

| framework [54] | | | | | | | |
|---|---|---|---|---|---|---|---|
| **R-CNN, VGG and RNN** [2] | 64.00 | 44.57 | 24.87 | 12.54 | - | - | - |
| **hybrid end-to-end CNN-LSTM networks** [55] | 69.6 | 57.0 | 46.9 | 40.0 | 34.2 | 163.2 | 68.5 |
| **Scene graph (CNN-RNN-SVM) framework** [6] | 67.6 | 49.3 | 35.5 | 26.1 | 22.3 | 76 | - |
| **CNN+RHN** [57] | 72.3 | 55.3 | 41.3 | 30.6 | 25.2 | 98.9 | - |
| **end-to-end deep learning approach** [5] | 66.8 | 46.8 | 32.2 | 22.1 | - | 55.12 | 20.36 |
| **Stack-VS** [58] | 79.0 | 63.4 | 48.9 | 37.2 | 28.8 | 118.9 | 57.5 |
| **P3DAT-CAS** [59] | 75.3 | 59.0 | 45.0 | 34.1 | 27.0 | 109.1 | 55.4 |
| **M2 Transformer** [60] | 81.6 | 66.4 | 51.8 | 39.7 | 29.4 | **129.3** | 59.2 |
| **Up-down** [61] | 80.2 | 64.1 | 49.1 | 36.9 | 27.6 | 117.9 | 57.1 |
| **Scene-based Factored Attention Module** [62] | 80.3 | 64.6 | 60.1 | 38.1 | 28.5 | 126.8 | 58.2 |

| ViT-GPT-2 (Our) | 92.18 | 86.15 | 81.17 | 74.19 | 45.36 | 81.37 | 59.5 |
|---|---|---|---|---|---|---|---|

The results are shown in Table 5.3 there is a clear difference between our model ViT-GPT-2 and other model's work on performance to perform image captioning tasks. However, after focusing on the semantic attention and visual attention of the image overall performance of our model improves and the scores of the various evaluation metrics improve which are higher than the other models that are mentioned in the literature review. Bleu-1, Bleu-2, Bleu-3, and Blue-4 are improved by 10.58, 20.45, 21.07, and 34.19 respectively. The other evaluation metrics like Meteor improve by 11.16 and the Rogue metric improves by 0.3. These results show that image captions generated by using the ViT-GPT-2 encoder-decoder framework are more reliable. Image captions overlap with the reference captions showing the performance of the caption generation the linguistic quality of the caption improves and all-important features of the images are extracted to generate a more accurate image caption. A graphical comparison of the performance of ViT-GPT-2 and another model on the MS COCO dataset is shown in Figure 5.3.
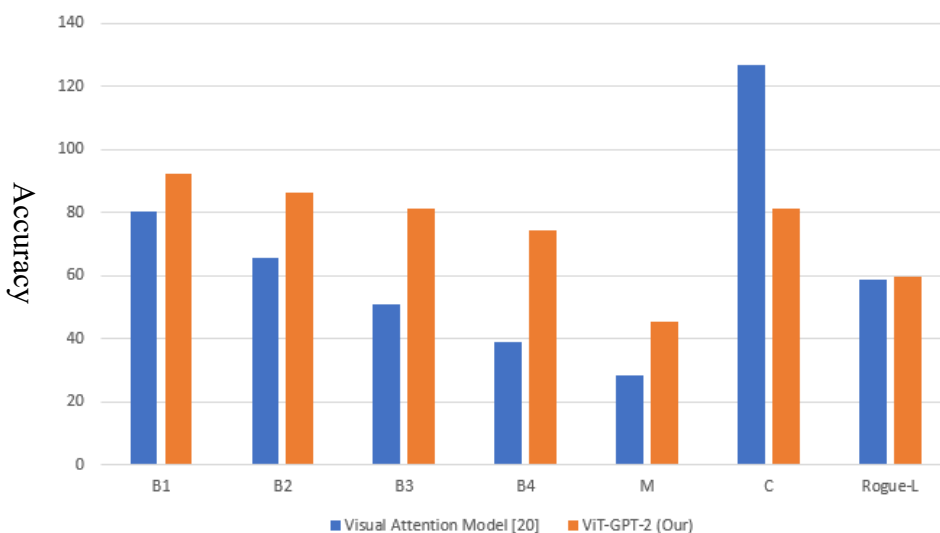


Evaluation Metrics

**Figure 5.3** Comparison with COCO Dataset

### 5.1.3 Comparison of Results with Base Paper on MS COCO Dataset

**Table 5.4** Result Comparison with Base Paper on MSCOCO Dataset

| Model Name | B1 | B2 | B3 | B4 | M | C | Rogue-L |
|---|---|---|---|---|---|---|---|
| **Visual Attention Model** [20] | 80.5 | 65.7 | 51.0 | 38.9 | 28.3 | 126.7 | 58.8 |
| **ViT-GPT-2 (Our)** | **92.18** | **86.15** | **81.17** | **74.19** | **45.36** | 81.37 | **59.5** |

The results are shown in Table 5.4 there is a clear difference between our model ViT-GPT-2 and base paper [20] on performance to perform image captioning tasks. Both models use same evaluation metrics like BLEU, METEOR, CIDEr and ROUGE. All of evaluation metrics improved except CIDEr graphical comparison is given in Figure 5.4.



Evaluation Metrics
**Figure 5.4** Comparison of Base Paper on COCO Dataset

### 5.1.4 Comparison of Our Results with Other Models Flicker 30k Dataset
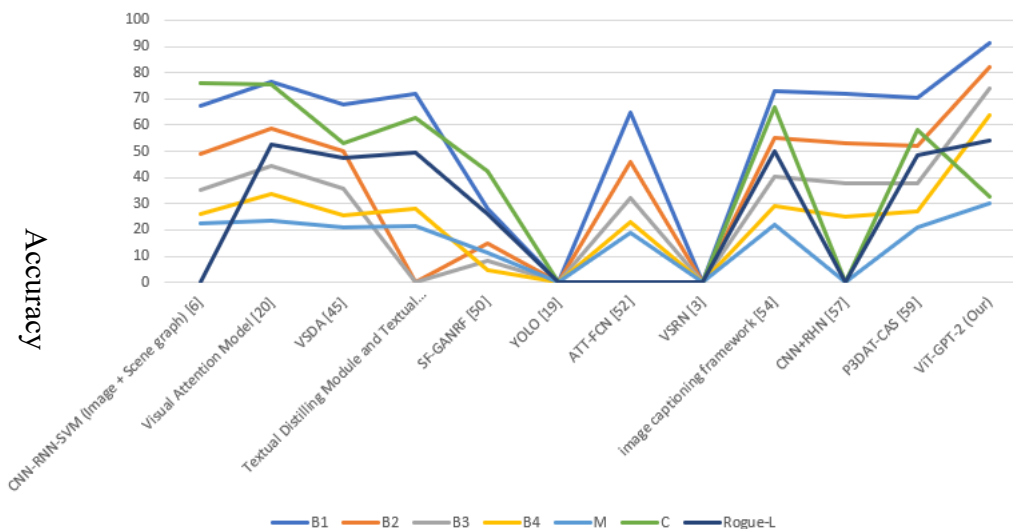
Table 5.4 presents the outcomes of our experiment conducted on the Flicker 30k dataset, showcasing a comparative analysis with other models. Our model places a significant emphasis on capturing the semantic knowledge of images during the image captioning process. It is noteworthy that distinct models exhibit varied behaviors, leading to diverse results in the conducted experiments.

**Table 5.5** Result Comparison with Other Models on Flicker 30k Dataset

| Model Name | B1 | B2 | B3 | B4 | M | C | Rogue-L |
|---|---|---|---|---|---|---|---|
| **CNN-RNN-SVM (Image + Scene graph)** [6] | 67.2 | 49.2 | 35.5 | 26.1 | 22.3 | 76 | - |
| **Visual Attention Model** [20] | 76.3 | 58.9 | 44.5 | 33.6 | 23.7 | 75.3 | 52.5 |
| **VSDA** [45] | 68.1 | 49.8 | 35.7 | 25.6 | 20.8 | 53.2 | 47.4 |
| **Textual Distilling Module and Textual Association Module** [47] | 71.8 | - | - | 27.9 | 21.6 | 62.7 | 49.3 |
| **SF-GAN$_{RF}$** [50] | 27.9 | 14.7 | 8.3 | 4.8 | 11.2 | 42.2 | 26.2 |
| **YOLO** [19] | -0.25 % | 0.45 % | -0.86 % | -1.63 % | 3.82 % | 12.63 % | 3.76 % |
| **ATT-FCN** [52] | 64.7 | 46.0 | 32.4 | 23.0 | 18.9 | - | - |
| **VSRN** [3] | 71.3 R@1 | 90.6 R@5 | 96.0 R@10 | - | - | - | - |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **image captioning framework** [54] | 73.1 | 55.1 | 40.1 | 29.0 | 22.0 | 66.8 | 50.1 |
| **CNN+RHN** [57] | 72.0 | 53.0 | 38.0 | 25.0 | - | - | - |
| **P3DAT-CAS** [59] | 70.3 | 52.2 | 37.8 | 27.1 | 21.1 | 58.3 | 48.4 |
| **ViT-GPT-2 (Our)** | **91.32** | **82.20** | **74.03** | **63.62** | **30.1** | 32.89 | **54.01** |

The results are shown in Table 5.4 these are the results calculated on the flicker 30k dataset and there is a clear difference between our model and other model's work performance in performing image captioning tasks. However, after focusing on the semantic attention and visual attention of the image overall performance of our model improves and the scores of the various evaluation metrics improve which are higher than the other models that are mentioned in the literature review. Bleu-1, Bleu-2, Bleu-3, and Blue-4 are improved by 15.02, 23.6, 29.53, and 30.02 respectively. The other evaluation metrics like Meteor improve by 6.4 and the Rogue metric improves by 1.51. A graphical comparison of the performance of ViT-GPT-2 and other models on the Flicker 30k dataset is shown in Figure 5.5.
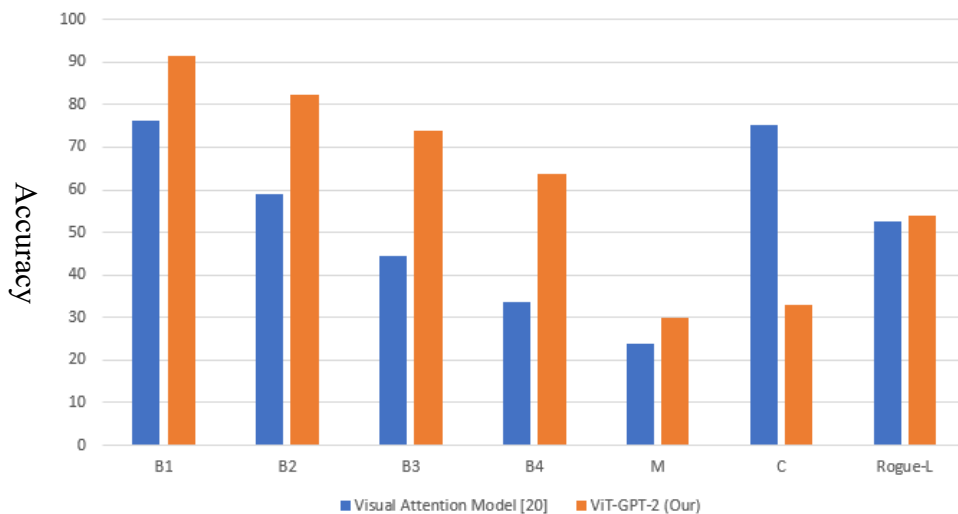
**Figure 5.5** Results comparisons on Flicker 30k Dataset

## 5.1.5 Comparison of Results with Base Paper on Flicker 30k Dataset

**Table 5.6** Result Comparison with Base Paper on Flicker 30k Dataset

| Model Name | B1 | B2 | B3 | B4 | M | C | Rogue-L |
|---|---|---|---|---|---|---|---|
| **Visual Attention Model** [20] | 76.3 | 58.9 | 44.5 | 33.6 | 23.7 | 75.3 | 52.5 |
| **ViT-GPT-2 (Our)** | **91.32** | **82.20** | **74.03** | **63.62** | **30.1** | 32.89 | **54.01** |

The results are shown in Table 5.6 there is a clear difference between our model ViT-GPT-2 and base paper [20] on performance to perform image captioning tasks. Both models use same evaluation metrics like BLEU, METEOR, CIDEr and ROUGE. All of evaluation metrics improved except CIDEr graphical comparison is given in Figure 5.6.

Evaluation Metrics
**Figure 5.6** Comparison of Base Paper on Flicker 30k Dataset

## 5.1.6 Qualitative Analysis

In this research, our main focus is on the improvement of performance of generating image caption. Lack of accuracy due to the lack of semantic knowledge while generating image captions. For the qualitative analysis of my research, images from the dataset were taken and captions by using our model and comparing these captions with the other model's generated captions with the grounded captions that are available in the dataset. All the images and captions and their comparison are given in Table 4.5. Our model ViT-GPT-2 improves the performance of the image caption and focuses on the semantics of the image both ViT and GPT-2 focus on the semantics of the image as well as the semantics of the text in this procedure. Detailed examples of the ViT-GPT-2 and other model captions are discussed in Table 5.8 (a, b, c, and d). In generated captions following points are discussed below:

- Captions generated through ViT-GPT-2 are focused on the linguistics of the sentences.
- ViT-GPT-2 detects all the objects in the image and creates a more relevant caption to the image.

- ViT-GPT-2 detects all the objects in the image and detects the time frame of the scene in the images.
- ViT-GPT-2 focuses on the scene understanding in the image and focuses on the position of objects while generating captions.

**Table 5.7** Captions comparison with other models on the MSCOCO dataset

| Image | Our Model Caption | Other Model Caption |
|---|---|---|
| | A bicycle parked in a grassy area next to a body of water | A Bike parked next to a bridge over a body of water |
| | A man sitting on a couch with a cat on his lap with laptop | A man sitting on a couch with a laptop computer |
| | a man riding a wave on top of a surfboard | a man riding a wave on top of a surfboard |

|  | a row of blue and yellow umbrellas on the side of a building | A group of color full umbrellas in front of building |
| --- | --- | --- |
|  | A man and a woman sitting on a sidewalk next to each other using mobile | A couple of women sitting on a bench looking at their mobile phone |
|  | a man and two women sitting at a table with wine glasses | A group of people sitting at a table with wine glasses |
|  | a dog sitting on a bed next to a pile of clothes | a dog sitting on a bed next to a group of clothes |

| | | |
|---|---|---|
|  | A city street with a traffic light and street signs | A group of cars on a city street with traffic signal |
|  | a man holding a frisbee in front of a building | A man holding a yellow frisbee in front of building |
|  | a woman standing in a living room holding a wii remote | A women is playing a video game in living room |
|  | a large clock tower towering over a city at night | A view of city street with a clock tower |

| | | |
|---|---|---|
|  | a woman sitting at a table with a lit candle | A group of people sitting on table with cake |
|  | a man holding a dog on top of a boat | A man sitting on a boat with a dog |
|  | a small dog sitting on top of a toilet seat | A dog sitting on top of a toilet seat. |
|  | a man riding a horse on top of a lush green hillside | A man riding a horse down a hill. |

Qualitative analysis of image captions generated by ViT-GPT-2, compared with captions from other models, involves a examination of the linguistic quality, coherence,

relevance, and creativity of the generated text. ViT-GPT-2's unique combination of Vision Transformer (ViT) for image understanding and GPT-2 for text generation allows for a comprehensive analysis of both visual and textual elements. Qualitative analysis provides valuable insights into the model's proficiency in generating human-like and contextually relevant image descriptions, enabling researchers to assess the advancements and limitations of ViT-GPT-2 concerning other state-of-the-art captioning models.

Table 5.8 shows that images and captions were generated by using the ViT-GPT-2 model and the other image captioning model. Here A qualitative analysis is performed on the image captions generated by ViT-GPT-2 and the other model and show which image captioning model generates more accurate image captions and focuses on the semantics of the image white generating image captions. Table 5.8 a, b, c, and d shows the comparison of the image captions and their qualitative analysis below that table.

**Table 5.8** Sample Images for Image Caption

Table 5.8 (a)

| Image | ViT-GPT-2 Caption | Other Model Caption |
|---|---|---|
|  | A women is sitting at a table with lit candle | A group of people sitting on table with cake |

**Accurate Object Detection:** Table 5.8 (a) shows the comparison of the image caption of ViT-GPT-2 and other model-generated captions. ViT-GPT-2 detects all the objects in the

image accurately like women, table and lit candle while other model did not detect objects accurately it detected cake instead of lit candle and detect group of people instead of women. Clarity in the object detection is due to the self-attention layer of transformer that focuses on the semantic of images. Caption generated by ViT-GPT-2 is more accurate as compare to the other model generated caption.

Table 5.8 (b)

| Image | ViT-GPT-2 Caption | Other Model Caption |
|---|---|---|
|  | A dog sitting on bed next to a pile of clothes | A dog is sitting on a bed next to a group of clothes |

**Improving Linguistic Quality:** Table 5.8 (b) shows the comparison of the image caption of ViT-GPT-2 and other model generated captions. ViT-GPT-2 detected all the objects and predicted words for the caption more accurately linguistically while other models did not focus on the linguistics of the sentence like ViT-GPT-2 generated the word 'pile of clothes' and other model generated 'group of clothes'. Caption generated by ViT-GPT-2 is accurate linguistically as compared to the other model generated caption because GPT-2 is trained on multiple books and many textual data therefore using GPT-2 linguistics of the captions improves. Improvement in the linguistic comes due to the self-attention layer of the GPT-2 that focus on the semantic of text while generating caption for image.

Table 5.8 (c)

| Image | ViT-GPT-2 Caption | Other Model Caption |
|---|---|---|
|  | A bicycle parked in a grassy area next to a body of water | A bike parked next to bridge over a body of water. |

**Object Detection With Object Position:** Table 5.8 (c) shows the comparison of the image caption of ViT-GPT-2 and other model generated captions. ViT-GPT-2 detected all the objects and understood the scene in the image accurately while other model did not detect objects accurately and also did not describe the scene accurately. ViT-GPT-2 detects bicycle, grass, and water surfaces and generates a caption that show the relationship between these objects. Other model detected bike instead of bicycle and missed some object and did not relate these objects' positions while generating an image caption. Image caption generated by ViT-GPT-2 are accurate and focus on scene understanding and relationships between the objects. Accurate object detection with their positions in the due to the semantic of the image and text that focused by the both the transformer models ViT and GPT-2.

Table 5.8 (d)

| Image | ViT-GPT-2 Caption | Other Model Caption |
|-------|-------------------|---------------------|
|  | A large clock tower towering over a city at night | A view of city street with a clock tower |

**Object Detection with Time Frame and Properties:** Table 5.8 (d) shows the comparison of the image caption of ViT-GPT-2 and other model generated captions. ViT-GPT-2 detects all the objects in the image detect the time frame of the scene that night time in the image and generate a caption sentence broadly like large tower and towering on the city. While other model did not detect the time frame and did not generate sentence broad. Caption generated by ViT-GPT-2 is more accurate as compare to the other model generated caption.

## 5.2  Discussion

All results are evaluated on the evaluation metrics that are used for the natural language processing such as BLEU, CIDEr, ROUGE and METEOR. BLEU scores measure the overlap of generated caption and reference caption and it indicates the similarity in the generated caption with reference caption. While METEOR metric evaluates the quality of the generated caption by focusing linguistic quality of the generated caption with the reference caption. It focusses not only exact words but it also focusses on the synonyms of the words. CIDEr considers the diversity and richness of vocabulary and concepts used in

the captions. ROUGE scores measure the longest overlap sentence of the generated caption and the reference caption. All metrics that are used for evaluations of the captions are improved except CIDEr.

All the scores improved due to the capability of our model that focus on the semantics of the image and text as well by using the self-attention layer of the transformer. By focusing on the semantics of the image and text following factors improved in this research

- Object Detection
- Object Detection with their relationships and sequence
- Object Detection with time frame and properties of objects e.g. day, night or long tower etc.
- Linguistic quality like pile of clothes rather than group of clothes

In image captioning, scenarios can arise where precision is high while recall is low. This imbalance typically occurs when machine-generated captions are exceptionally accurate but may miss describing some elements present in the image. For instance, consider an image showing a beach scene with various objects such as umbrellas, people, and waves. If a machine-generated caption accurately describes the people and umbrellas but omits mentioning the waves, precision remains high due to the correctness of the mentioned elements. However, since the caption fails to encompass all relevant details present in the image, recall suffers. In essence, high precision with low recall signifies that the generated captions are precise and accurate concerning the described elements but might not fully capture the entirety of the image's content.

In image captioning, it's possible for BLEU to yield higher scores compared to METEOR due to differences in their calculation methods and focuses. BLEU primarily emphasizes n-gram overlap between machine-generated and reference captions, rewarding exact matches and penalizing deviations from reference captions. As a result, if a machine-generated caption closely matches reference captions in terms of specific n-grams, BLEU tends to assign a high score. For example, if a reference caption states "a dog running in the park," and the machine-generated caption reads "a dog running in the park," BLEU would assign a high score due to the exact match.

On the other hand, METEOR considers additional factors such as synonymy, stemming, and word order in its evaluation, aiming to capture a broader range of linguistic phenomena beyond exact word matches. Therefore, METEOR might assign lower scores if the machine-generated caption diverges from reference captions in terms of word choice, word order, or the inclusion of synonyms and related words. For instance, if the reference caption mentions "a canine jogging in the park," and the machine-generated caption reads "a dog running in the park," METEOR may assign a lower score due to differences in word choice (i.e., "canine" vs. "dog").

## 5.3 Summary

In this chapter, A complete quantitative and qualitative result analysis of the research was discussed. At the end of this chapter comparison of results with other researcher that work in the domain of 'image captioning'. The captions generated by ViT-GPT-2 were compared with the captions generated by the other researchers using other models.

# CHAPTER 6

# CONCLUSION AND FUTURE DIRECTIONS

## 6.1  Conclusion

In this study, the possibility for collaboration when the Vision Transformer (ViT) and Generative Pre-trained Transformer 2 (GPT-2) were used together to tackle the difficult task of 'Image Captioning'. A novel framework developed by me that smoothly combines visual and textual information by utilizing ViT's capacity to extract fine-grained visual elements and GPT-2's expertise in language creation. Through empirical analyses, it demonstrates by me that how well this method works at producing cohesive, contextually appropriate captions that accurately represent the content of a variety of photos.

My experiments showed that the ViT-GPT-2 combination takes advantage of both models to handle the challenges of image captioning. ViT accurately catches the subtleties of the visual content of images, while GPT-2 skillfully transforms these visual elements into concise descriptions. By merging vision and language models, this multimodal synergy advances the state of the art in image captioning. The results are encouraging.

According to my research, the Vision Transformer successfully collects complex visual details at many scales, enabling it to comprehend both the local and global components of the images. For GPT-2 to be able to generate captions, this feature extraction technique is essential in giving pertinent visual context. In contrast, GPT-2 demonstrates

its skill at language generation by skillfully fusing logical and fluid textual descriptions with the visual indications offered by ViT.

An experiment performed on the MS COCO dataset and results evaluated on the evaluation metrics like BLEU-1, BLEU-2, BLEU-3, BLEU-4, METEOR, CIDEr, and Rogue and results come 91.21, 86.15, 81.17, 74.19, 45.36,81.37 and 559.5 respectively. Improvements to Blue-1, Blue-2, Blue-3, and Blue-4 are 10.58, 20.45, 21.07, and 34.19, respectively. Other evaluation criteria such as Meteor and Rogue improve by 11.16 and 0.3, respectively.

## 6.2   Future Directions

For the future directions in the domain of image captioning to enhance the performance of image captioning. Investigate methods to more efficiently combine the visual and textual embeddings, allowing the model to provide captions that better capture complex links between the visual and linguistic features. Increase the model's capacity to identify and characterize the finer points of an object, allowing for more detailed and specialized captions. The area of image captioning employing the ViT-GPT-2 architecture can advance by exploring these future paths to provide more precise, inventive, and contextually relevant captions that satisfy a variety of applications and user needs.

# References

[1]    T. Ghandi, H. Pourreza, and H. Mahyar, "Deep learning approaches on image captioning: A review," *ACM Comput Surv*, vol. 56, no. 3, pp. 1–39, 2023.

[2]    X. Zhang, X. Song, X. Lv, S. Jiang, Q. Ye, and J. Jiao, "Rich image description based on regions," in *Proceedings of the 23rd ACM international conference on Multimedia*, 2015, pp. 1315–1318.

[3]    K. Li, Y. Zhang, K. Li, Y. Li, and Y. Fu, "Visual semantic reasoning for image-text matching," in *Proceedings of the IEEE/CVF International conference on computer vision*, 2019, pp. 4654–4662.

[4]    zakir Hossain, F. Sohel, F. M. Shiratuddin, and H. Laga, "A Comprehensive Survey of Deep Learning for Image Captioning," *ACM*, vol. 51, no. 6, pp. 1–36, Feb. 2109.

[5]    B. Wang, C. Wang, Q. Zhang, Y. Su, Y. Wang, and Y. Xu, "Cross-lingual image caption generation based on visual attention model," *IEEE Access*, vol. 8, pp. 104543–104554, 2020.

[6]    L. Gao, B. Wang, and W. Wang, "Image captioning with scene-graph based semantic concepts," in *Proceedings of the 2018 10th international conference on machine learning and computing*, 2018, pp. 225–229.

[7]    J. Wei, Z. Li, J. Zhu, and H. Ma, "Enhance understanding and reasoning ability for image captioning," *Applied Intelligence*, vol. 53, no. 3, pp. 2706–2722, 2023.

[8]    M. Hossain and others, "Deep learning techniques for image captioning," Murdoch University, 2020.

[9]    A. Jamil *et al.*, "Deep Learning Approaches for Image Captioning: Opportunities, Challenges and Future Potential," *IEEE Access*, 2024.

[10]   H. Sharma, M. Agrahari, S. Kumar Singh, M. Firoj, and R. Kumar Mishra, "Image Captioning: A Comprehensive Survey," in *2020 International Conference on Power*

*Electronics & IoT Applications in Renewable Energy and its Control (PARC)*, Mathura: IEEE, May 2020.

[11]   P. Zhao *et al.*, "Retrieval-Augmented Generation for AI-Generated Content: A Survey," *arXiv preprint arXiv:2402.19473*, 2024.

[12]   A. Siddharthan, "Ehud Reiter and Robert Dale. Building Natural Language Generation Systems. Cambridge University Press, 2000. 64.95/£ 37.50 (Hardback). 234 pages," *Nat Lang Eng*, vol. 7, no. 3, pp. 271–274, 2001.

[13]   E. B. Goldstein, *Cognitive psychology: Connecting mind, research and everyday experience*. Cengage Learning, 2014.

[14]   K. Cherry, "Cognition in psychology." Accessed: Jun. 13, 2023. [Online]. Available: https://www.verywellmind.com/what-is-cognition-2794982#toc-uses-of-cognition

[15]   F. Donnarumma, H. Dindo, and G. Pezzulo, "Sensorimotor coarticulation in the execution and recognition of intentional actions," *Front Psychol*, vol. 8, p. 237, 2017.

[16]   L. Song, F. Li, Y. Wang, Y. Liu, Y. Wang, and S. Xiang, "Image captioning: Semantic selection unit with stacked residual attention," *Image Vis Comput*, p. 104965, 2024.

[17]   S. Russell, *Artificial Intelligence: A Modern Approach, eBook, Global Edition*. Pearson Education, Limited, 2016.

[18]   J. Bach, *Principles of synthetic intelligence PSI: an architecture of motivated cognition*, vol. 4. Oxford Cognitive Models and Ar, 2009.

[19]   M. A. Al-Malla, A. Jafar, and N. Ghneim, "Image captioning model using attention and object features to mimic human image understanding," *J Big Data*, vol. 9, no. 1, pp. 1–16, 2022.

[20]   H. Wei, Z. Li, F. Huang, C. Zhang, H. Ma, and Z. Shi, "Integrating scene semantic knowledge into image captioning," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 17, no. 2, pp. 1–22, 2021.

[21]   S. Dargan, M. Kumar, M. R. Ayyagari, and G. Kumar, "A survey of deep learning and its applications: a new paradigm to machine learning," *Archives of Computational Methods in Engineering*, vol. 27, pp. 1071–1092, 2020.

[22]   S. Saha, "A comprehensive guide to convolutional neural networks—the ELI5 way," *Towards data science*, vol. 15, p. 15, 2018.

[23]   K. Pathak, M. Pavthawala, N. Patel, D. Malek, V. Shah, and B. Vaidya, "Classification of brain tumor using convolutional neural network," in *2019 3rd International conference on Electronics, Communication and Aerospace Technology (ICECA)*, 2019, pp. 128–132.

[24]   R. M. Oruganti, *Image description using deep neural networks*. Rochester Institute of Technology, 2016.

[25]   A. Graves, M. Liwicki, H. Bunke, J. Schmidhuber, and S. Fernández, "Unconstrained on-line handwriting recognition with recurrent neural networks," *Adv Neural Inf Process Syst*, vol. 20, 2007.

[26]   S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput*, vol. 9, no. 8, pp. 1735–1780, 1997.

[27]   H. Wang, L. Jiang, Q. Huang, and L. Zhao, "Image Captioning with Object Detection and Localization," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 180–196.

[28]   A. Vaswani *et al.*, "Attention is All You Need," in *Advances in Neural Information Processing Systems (NIPS)*, 2017.

[29]   V. Mnih, N. Heess, A. Graves, and others, "Recurrent models of visual attention," *Adv Neural Inf Process Syst*, vol. 27, 2014.

[30]   O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3156–3164.

[31]  J. Wu, H. Hu, and L. Yang, "Pseudo-3D attention transfer network with content-aware strategy for image captioning," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 15, no. 3, pp. 1–19, 2019.

[32]  A. Vaswani *et al.*, "Attention is All You Need," *CoRR*, vol. abs/1706.03762, 2017.

[33]  K. Luu, X. Wu, R. Koncel-Kedziorski, K. Lo, I. Cachola, and N. Smith, "Explaining Relationships Between Scientific Documents," in *Proceedings of the Conference on Knowledge Discovery and Data Mining (KDD)*, 2023.

[34]  J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.

[35]  W. A. Nanal, "Image Captioning Using Transformer Architecture," Purdue University, 2022.

[36]  R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.

[37]  J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.

[38]  A. Vaswani *et al.*, "Attention is all you need," *Adv Neural Inf Process Syst*, vol. 30, 2017.

[39]  A. Dosovitskiy *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[40]  A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.

[41] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3128–3137.

[42] K. Cho *et al.*, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.

[43] I. Sutskever, O. Vinyals, and Q. V Le, "Sequence to sequence learning with neural networks," *Adv Neural Inf Process Syst*, vol. 27, 2014.

[44] J. Mao, W. Xu, Y. Yang, J. Wang, Z. Huang, and A. Yuille, "Deep captioning with multimodal recurrent neural networks (m-rnn)," *arXiv preprint arXiv:1412.6632*, 2014.

[45] C. He and H. Hu, "Image captioning with visual-semantic double attention," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 15, no. 1, pp. 1–16, 2019.

[46] B. Zhang *et al.*, "Structural semantic adversarial active learning for image captioning," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 1112–1121.

[47] F. Liu, X. Wu, S. Ge, X. Zhang, W. Fan, and Y. Zou, "Bridging the gap between vision and language domains for improved image captioning," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 4153–4161.

[48] S.-H. Han and H.-J. Choi, "Domain-specific image caption generator with semantic ontology," in *2020 IEEE International Conference on Big Data and Smart Computing (BigComp)*, 2020, pp. 526–530.

[49] Y. Yang, "Image-Caption Pair Replacement Algorithm towards Semi-supervised Novel Object Captioning," in *2022 7th International Conference on Intelligent Computing and Signal Processing (ICSP)*, 2022, pp. 266–273.

[50] J. Feng and J. Zhao, "Improving Caption Consistency to Image with Semantic Filter by Adversarial Training," in *2021 IEEE International Conference on Mechatronics and Automation (ICMA)*, 2021, pp. 269–274.

[51] Z. Zhang, Q. Wu, Y. Wang, and F. Chen, "High-quality image captioning with fine-grained and semantic-guided visual attention," *IEEE Trans Multimedia*, vol. 21, no. 7, pp. 1681–1693, 2018.

[52] Q. You, H. Jin, Z. Wang, C. Fang, and J. Luo, "Image captioning with semantic attention," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4651–4659.

[53] K. Nguyen, S. Tripathi, B. Du, T. Guha, and T. Q. Nguyen, "In defense of scene graphs for image captioning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1407–1416.

[54] W. Cai and Q. Liu, "Image captioning with semantic-enhanced features and extremely hard negative examples," *Neurocomputing*, vol. 413, pp. 31–40, 2020.

[55] A. Nguyen, Q. D. Tran, T.-T. Do, I. Reid, D. G. Caldwell, and N. G. Tsagarakis, "Object captioning and retrieval with natural language," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, p. 0.

[56] L. Gao, B. Wang, and W. Wang, "Image captioning with scene-graph based semantic concepts," in *Proceedings of the 2018 10th international conference on machine learning and computing*, 2018, pp. 225–229.

[57] J. Gu, G. Wang, J. Cai, and T. Chen, "An empirical study of language cnn for image captioning," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1222–1231.

[58] L. Cheng, W. Wei, X. Mao, Y. Liu, and C. Miao, "Stack-VS: Stacked visual-semantic attention for image caption generation," *IEEE Access*, vol. 8, pp. 154953–154965, 2020.

[59] J. Wu, H. Hu, and L. Yang, "Pseudo-3D attention transfer network with content-aware strategy for image captioning," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 15, no. 3, pp. 1–19, 2019.

[60] M. Cornia, M. Stefanini, L. Baraldi, and R. Cucchiara, "Meshed-memory transformer for image captioning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 10578–10587.

[61] P. Anderson *et al.*, "Bottom-up and top-down attention for image captioning and visual question answering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6077–6086.

[62] C. Shen, R. Ji, F. Chen, X. Sun, and X. Li, "Scene-based Factored Attention for Image Captioning," *arXiv preprint arXiv:1908.02632*, 2019.

[63] T.-W. Wu, J.-H. Huang, J. Lin, and M. Worring, "Expert-defined keywords improve interpretability of retinal image captioning," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 1859–1868.

[64] N. Wang, J. Xie, J. Wu, M. Jia, and L. Li, "Controllable image captioning via prompting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023, pp. 2617–2625.

[65] J. Luo *et al.*, "Semantic-conditional diffusion networks for image captioning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 23359–23368.

[66] Y. Gong, L. Wang, M. Hodosh, J. Hockenmaier, and S. Lazebnik, "Improving image-sentence embeddings using large weakly annotated photo collections," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part IV 13*, 2014, pp. 529–545.

[67] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik, "Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2641–2649.

[68] T.-Y. Lin *et al.*, "Microsoft coco: Common objects in context," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, 2014, pp. 740–755.

[69] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002, pp. 311–318.

[70] C. Callison-Burch, M. Osborne, and P. Koehn, "Re-evaluating the role of BLEU in machine translation research," in *11th conference of the european chapter of the association for computational linguistics*, 2006, pp. 249–256.

[71] E. Denoual and Y. Lepage, "BLEU in characters: towards automatic MT evaluation in languages without word delimiters," in *Companion Volume to the Proceedings of Conference including Posters/Demos and tutorial abstracts*, 2005.

[72] S. Banerjee and A. Lavie, "METEOR: An automatic metric for MT evaluation with improved correlation with human judgments," in *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005, pp. 65–72.

[73] R. Vedantam, C. Lawrence Zitnick, and D. Parikh, "Cider: Consensus-based image description evaluation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 4566–4575.

[74] S. Robertson, "Understanding inverse document frequency: on theoretical arguments for IDF," *Journal of documentation*, vol. 60, no. 5, pp. 503–520, 2004.

[75] P. Anderson, B. Fernando, M. Johnson, and S. Gould, "Spice: Semantic propositional image caption evaluation," in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part V 14*, 2016, pp. 382–398.

[76] J. Johnson *et al.*, "Image retrieval using scene graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3668–3678.

[77] S. Schuster, R. Krishna, A. Chang, L. Fei-Fei, and C. D. Manning, "Generating semantically precise scene graphs from textual descriptions for improved image retrieval," in *Proceedings of the fourth workshop on vision and language*, 2015, pp. 70–80.

[78] C.-Y. Lin, "Rouge: A package for automatic evaluation of summaries," in *Text summarization branches out*, 2004, pp. 74–81.