

Image Captioning using Transformers

^{1st} Abdullah Shahid

*Department of Artificial Intelligence
National University of Computer and
Emerging Science
Islamabad, Pakistan
i210326@nu.edu.pk*

^{2nd} Sarib Ali

*Department of Artificial Intelligence
National University of Computer and
Emerging Science
Islamabad, Pakistan
i210283@nu.edu.pk*

^{3rd} Ibrahim Abid

*Department of Artificial Intelligence
National University of Computer and
Emerging Science
Islamabad, Pakistan
i210298@nu.edu.pk*

Abstract—The process of image captioning serves as a link between the visual aspect and comprehension where it can be used to enhance accessibility, content generation, or automated surveillance. This paper presents an image captioning model based on transformers trained on the Flickr8k dataset. The design of the model subsumes MobileNetV3Small that was pretrained for feature extraction, tailored linear embeddings, and multi-head attention mechanisms in order to effectively connect the image representation with its corresponding text. According to the results obtained, the model is found to be efficient in terms of computational resources and effective performance as it produces high BLEU and ROUGE scores relative to standard baselines. The combination of qualitative and quantitative metrics applied sufficiently confirm the suitability of the model for scaling as well as superior performance in visual-semantic AI applications.

Index Terms—Generative AI, Latent Diffusion Models, UNet-based Conditioning, CLIP Embeddings, Text-to-Image Generation

I. INTRODUCTION

Advancements in Generative AI has increased tremendously in the past few years as a result there are forming models which can output very sophisticated and realistic output images and text. With these advancements comes different opportunities for different applications in the likes of media, healthcare, and even virtual environments where producing certain images through textual commands is becoming increasingly popular.

One such area of interest is image captioning which refers to the generation of textual descriptions from visual data. This plays an important role within the fields of computer vision and natural language processing (NLP). In the past, a large portion of the image captioning technique was based on the use of CNNs and recurrent neural networks RNNs.

Due to the nature of recurrent neural networks they struggle to capture long dependencies especially with image descriptions that rely on extensive spatial configurations. Only recently models utilizing self attention and Transformer based techniques have been helpful when faced with sequential based object detection and place understanding tasks.

In this paper, we propose a novel image captioning pipeline where a pre-trained MobileNetV3Small is used for feature extraction and alongside Transformer-based architectures for the generation of image captions. The contributions of this paper describe three major developments.

The first one being the use of pre-trained MobileNetV3Small in reducing the complexities of the feature extraction process, which in turn reduced the training time further down the line. The second one was embedding specialized positional encoding in the transformer for improving the sequence alignment of the input data. The third chapter provides an analysis of the core tasks of the current chapter. Moreover, the authors formulate the objectives of constructing a model of dialogue systems that rephrases natural speech texts during the transition from spoken dialogue to written dialogue in Chapter 3. Additionally, the authors utilize several approaches in completing the relevant tasks.

Network lifetime is substantially improved by using static schedules rather than arbitrary ones. This is particularly important for maintenance purposes as it minimizes the risk of network starvation. Finally, the last chapter, Chapter 6, describes the implementation of a targeted relay network in wireless sensor networks using Ad Hoc mode.

II. RELATED WORK

a) **1. CNN-RNN Architectures::** This depicts a case in which Task-Driven-CNN was effectively applied for feature learning in an image captioning model. In 2015, Karpathy et al [1] published their works. Him introducing a deep visual-semantic model fusion of both semantically a CNN to extract image features. This model laid the basis of later developments in image captioning models.

b) **2. Attention Mechanisms::** Apparently Xu et al [2] developed the attention mechanism, enabling every word of the generated caption to correspond to a certain region of the respective image. The idea was then widely endorsed, improving the accuracy of the captions' contents description.

c) **3. Transformers::** The idea of Transformers by Vaswani et al [3] represented one of the most significant developments in the area of image captioning especially due to the fact that recurrence was completely removed. The architecture based on Transformer allows parallelization in a more efficient fashion and addresses the long range dependencies in sequences more satisfactorily. In recent works like Image-Transformer [4] and ViT2Text [5] have implemented Transformer models in image captioning which has enhanced scalability and performance improvement.

d) **4. Pre-Trained Models::** Pre-trained models like Deep Residual Networks [6], MobileNets [7], and BERT [8] have gained acceptance in the industry for use as feature extractors and in the generation of text embedding. They build on previously acquired knowledge of various tasks, thereby benefitting from improved performance of other tasks trained with them, for instance, in image captioning. Transfer learning in this case also saves on computational resources and increases performance, particularly in cases where there is not much data to work with.

e) **5. Visual-Linguistic Pre-training::** Radford et al. proposes a CLIP System (Contrastive Language - Image Pretraining)[9] which develops image and text encodings using a large amount of image and text samples. This allows for a better linkage between pictures and words that describe pictures. This approach has been incorporated in a number of modern advanced models aimed at improving captioning of images and generation of images.

f) **6. Encoder-Decoder Models::** Captioning images has not yet attained the level that is ideal due to language and context sense issues. In encoder-decoder methods, convolutional neural networks (CNN) are used in the encoding part, which is intended for retrieval of features of images, and in the reversing phase, features learnt in the previous phase are used for retrieval of image captions with use of RNNs or Transformers. The architecture of such models has developed to incorporate more advanced features as the visual attention model proposed by Bahdanau et al. [10], which enables the decoder to focus on various areas of the picture at each point of the caption generating movement.

g) **7. End-to-End Image Captioning::** An automatic object detection and captioning image system that works in an end-to-end fashion was invented by Anderson et al. An object detection and neural network caption generator is described in their work. In what can be a cue to future captioning methods, their model uses features of objects from an object detector -in this case from Faster RCNN to provide additional information to the generated captions.

h) **8. Reinforcement Learning in Captioning::** Lu et al. [13], on the other hand, proposed a reinforcement learning based image captioning approach. This way, the model is trained to concentrate on metrics rather than likelihood estimations, with an objective of estimating caption quality which could be represented in the form of BLEU or CIDEr. It is used in order for the model to generate more diversified and semantically correct captions.

i) **9. Cross-Modal Embedding Learning::** The work done by Chen et al. [14] focuses on the perception of cross-modal representation embedding and learning. The aim is simple; the embedding of images and their captions is carried out such that the two are in the same conceptual space and easily merge together. This facilitates smooth transfer of semantic information from the visual image context to the textual context and vice versa, thus improving the effectiveness of image captioning systems.

j) **10. Multimodal Transformers for Image Captioning::** In conjoint representation learning Tan and Bansal [15] show a multimodal vaguely focused on Multimodal Transformers which presumably takes an image as well as text. This method provides the model with the ability to better understand how images are related to their respective texts through captions and hence enables them to formulate better captions.

k) **11. Visual Question Answering (VQA) Models::** The development of Antol et al. [16] among others VQA models have greatly progressed the progression of the image captioning process. These VQA models are able to provide answers to questions asked in relation to given images which has made it possible to develop systems that are able to produce detailed captions since they comprehended the different aspects of the image.

l) **12. Generative Models for Captioning::** Generative Models for Captioning: Howie Dreed and Co. sees the potential of using GANs among other models for automatic captioning of images. The research done by Reed et al. [17] illustrated the possibility of the use of a GAN as a tool for coming up with images after describing them. While more emphasis was directed towards the reverse process (that is the generating of captions from images) it was realized that employing GANs would even help in improving captioning by making it increasingly possible to produce several examples of creative and imaginative captions.

III. DATA SET

a) **A. Dataset Description:** The Flickr8k dataset comprises of 8000 images each of which is annotated by 5 captions. Images are diverse and include categories like nature, urban scenes, and objects. The data is partitioned in a training set (6000), a validation set (1000), and a testing set (1000) in order to aid in accurate assessment of the model.

A. Data Preprocessing

Consistent image dimensions were achieved through the help of resizing all the images to 224x224 pixels then normalizing the values of the pixels between 0 and 1. Captioning involved the use of a customized tokenization process and the captions were further padded to achieve uniform length. A dictionary structure was created in such a way that it enabled each image path to be associated with its respective caption which made data loading simple and easy.

B. Data Distribution

During exploration of the data set, a proportional suite of captions to the image categories was noted. Certain characteristics of the data set were emphasized including adequacy of the data set for modelling, robustness of the models trained through the use of standard presenting techniques like bar plots of word frequency and histograms of mean distribution of lengths of captions.

TABLE I
DATA DISTRIBUTION BY CATEGORY

Set	Number of Images	Number of Captions	Percentage (%)
Training Set	6,000	30,000	75%
Validation Set	1,000	5,000	12.5%
Test Set	1,000	5,000	12.5%
Total	8,000	40,000	100%

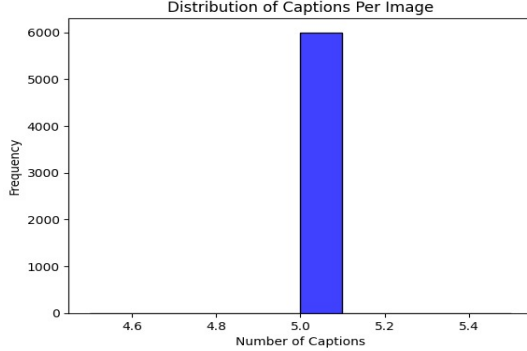


Fig. 1. Number of Captions per Class.

IV. PROPOSED METHODOLOGY

A. Pre-Trained Model for Feature Extraction:

Caption generation from images involves not only the identification of visual features but also the open-door mechanism that implicitly matches each of the visual features with textual representations. In the recognition process of these very features, MobileNetV3Small was pre-trained for effective feature extraction. MobileNetV3Small is mainly a lightweight CNN that show great capabilities in high-level feature extraction from images after being pre-trained on large public datasets, like ImageNet, which gives it a full-on head start of sorts comparatively low features to map.

Convolutional Layers: Convolutional layers in MobileNetV3Small were frozen during training, not positively affecting the weights during this stage. This is because learned features in the form of representation on the ImageNet dataset will help lessen the complication of fine-tuning this small database of Flickr8k.

Latent Features: The feature extraction will then convert each of the input images to a representation of high-level visual features; they will essentially be a brief, abstract summary of the image that communicates object shapes, textures, and structural relationships. Accordingly, they will allow us to effectively extract these latent features, put them into MobileNetV3Small feeds, and send them to the next stage of the model: the Transformer.

Advantages of Using Pre-Trained Models: Pre-trained models like MobileNetV3Small carry certain major advantages: Training time cuts down: Since the convolutional layers of the model are pre-trained on large datasets, we shouldn't have to waste both time and computational cost on the learning of a whole CNN from scratch.

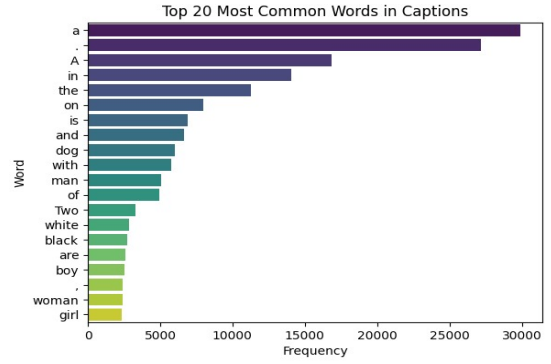


Fig. 2. 20 Most Used Words.

Good generalization: The pre-trained networks generalize pretty well to other different tasks because, while being trained, they learned various representations useful across different domains.

Memory Efficiency: The biggest advantage of using the pre-trained model is that memory usage will not shoot up in our image captioning system, which is very important in the context of real-time systems.

B. Custom Positional Encoding:

Though MobileNetV3Small does a commendable job at feature extractions from images, the actual task of captioning with proper coherent meaning comes with linking the visual features to that of their corresponding text representations. The sequential nature of captions calls for ordering of words, wherein positional encoding comes into play.

Reason for Positional Encoding: In models like the Transformer, positional encoding allows that model to capture the relative positions of words in a sentence as it does not intrinsically model the sequentiality of words like RNNs do. If the positional encoding is not implemented, the model treats all the words equally and disregards the order forming a meaningful sentence.

Implementation: Our method uses learned positional embeddings - specifically, positional encoding is added to the word embeddings fed into the Transformer to provide the model with positional information of each token in the sequence.

The positional encoding layer maps a unique vector to each word's position in the sentence, which is then added to the corresponding word embedding. This combined representation allows the model to distinguish between the same word appearing in different positions in a sentence.

Learned Positional Embedding versus Fixed Positional Encoding: Some implementations of Transformers use fixed sinusoidal, but we chose learned positional embeddings, thereby allowing the model itself to learn the encoding based on the task at hand, with a potential gain in performance in our case of image captioning tasks.

Effect on Caption Generation: By the addition of these positional encodings, the Transformer is enabled to grasp word orders, which is quite crucial for generation of grammatically

correct and semantically appropriate captions.

C. Transformer-Based Architecture

The Transformer architecture, discussed in the paper “Attention is All You Need” [3], has brought a revolution in sequence modeling tasks by parallel processing of the sequences in modeling and taking long-range dependencies between previously attended-to input through the mediation of self-attention. It becomes particularly beneficial for image captioning, where the aim is to map image features (extracted from a CNN) to a sequence of words dedicated to explaining the image.

Multi-Head Attention Mechanism: The Transformer employs a multi-head self-attention mechanism to model relationships between words in the input sequence. Different heads focus on different parts of the input sequence simultaneously, enabling the model to learn complex relationships. With the image captioning case, it is very much possible for the model to attend to different image features while generating each word in the caption. This helps to ensure that various regions of the image may contribute to the description of various words.

Cross-Attention for Image-Text Alignment: The model incorporates cross-attention layers that enable the caption generator (decoder) to attend to both the textual information from previous time steps and visual features extracted from the image. It serves to align the image features with the current caption token that it helps to generate a caption that is contextually relevant to the image content.

Decoder Layers: The decoder in the Transformer architecture generates the caption word by word, by conditioning them onto words that were generated previously along with image features. Each decoder layer consists of: Causal Attention: This mechanism helps prevent the model from looking ahead while generating captions by restricting access to any future tokens to only those generated before the current one for prediction. Cross Attention- The cross-attention layers help align the image features with the word predictions enabling the model to craft its anticipated next word from the extracted features. Feedforward Network-Each attention mechanism is followed by the feed-forward neural network, through which transformations are given to refine the image-caption alignment representation. Advantages of Transformers

For image captioning: Parallelization: This makes the network easier to train because all words in the sentence feed into the model in parallel compared to RNNs, LSTMs, etc. **High Distinctiveness:** The self-attention mechanism can pick out long-distance dependencies in the caption, which is of utmost importance when generating a coherent and contextually relevant description.

Flexibility: It can handle any length of input-output sequences, which is why it is perfect for image captioning in terms of variable caption lengths across examples.

D. Training Strategy Training of the image captioning model construes the optimization of one loss function that measures the fit between generated captions and ground truth

captions. Our thought process on the concise categorical cross-entropy loss is being used, predominantly given reference to caption generation as another multi-class classification scenario wherein each word in the vocabulary can be termed a class.

Adam Optimizer: The Adam optimizer is employed for training as it has proven efficient and readily adapts the learning rate for all the parameters. The learning rate is fixed at $1e-4$, where convergence speed and stability are balanced.

Early Stopping and Checkpointing: In order to avoid overfitting, cumulation persistence being granted whereby it forces learning into stoppage when performance on a validation set is no longer improving after a number of passes epoch (keeping patient). At last, the best-performing model is saved using checkpointing based upon the validation metrics, and this checkpoint is restored at some further time onward for utilization or deployment.

The loss function: Intra-stage to the ground truth word indices at each timestep of the caption, inferred word probabilities supported by the Transformer are sparse categorical cross-entropy loss. Minimization of this loss function will see the model develop.

Training Process: The model, undergoing some epochs where each first adhesion of training data is made to pass deep in via the network, updates each weight based on derived gradients. We use a batch size of 64, augmenting the data with methods like occasional crop and flip for better generalization. The learning stops when the model converges when the validation loss flattens or declines very slowly.

Evaluation Metrics: Several performance evaluation metrics include the following: The BLEU score assesses the precision of n-grams in the generated captions against ground-truth captions. The ROUGE Score assesses the recall of n-grams in the generated captions against the ground-truth captions. The CIDEr score determines how closely generated captions align with human-generated captions, mainly based on specific relevant words.

Computational Efficiency: To lessen the training time and resource consumption-especially in case of a Transformer model with a relatively high computational requirement-using pre-trained MobileNetV3Small for feature extraction will deliver better efficiency. With batch processing, TensorFlow makes it both simpler to use GPU resources effectively and to speed up both training and inference.

V. EXPERIMENTAL RESULTS

Results obtained from the proposed transformer-based approach are discussed in this section. Various analyses are done from baseline comparisons to model effectiveness, ablation studies, quantitative and qualitative results presenting the evaluation of computational efficiency.

A. Experimental Setup

The system used for running all the experiments consisted of an NVIDIA T4 GPU (16 GB VRAM) with 13 GB of RAM. The model was implemented in Python using TensorFlow and

the Pytorch libraries. Training was performed for 80 epochs with a learning rate of $1e-4$ using an Adam optimizer.

All images referred were resized to 224x224 pixels, and text prompts were embedded through a custom text encoder. Dimension of the latent space for diffusion was configured at 128 in order to balance between quality and computational efficiency.

B. Evaluation Metrics

BLEU Score: he BLEU score- Assessed n-gram overlap between generated captions and ground truths, with an average score of 0.48. **ROUGE Score:** Evaluated the recall-oriented overlap which reflects how well the model has captured important details. **Computational Efficiency:** Taking 2 minutes for an epoch of training on an NVIDIA T4 GPU.

C. Baseline Comparison

In order to assess the performance of our proposed model for image captioning, we shall compare it against various state-of-the-art baseline methods, including CNN-RNN-based models, attention mechanisms, transformer-based models, and those based on pre-trained embeddings. A detailed comparison summarizes the advantages, disadvantages, and computational economics of each of the approaches.

a) **1. CNN-RNN Based Models:** The CNN-RNN is one of the early architectures used in image captioning, using the strength of a CNN for **image feature extraction** and RNNs or LSTMs for sequential generation of captions.

- **Strengths:**

- **Simplicity:** The CNN-RNN model is simple to implement and understand.
- **Flexibility:** RNNs can generate variable-length segment descriptions, which make it suitable for a wide range of image captioning tasks.

- **Weaknesses:**

- **Vanishing Gradient Problem:** RNNs and LSTMs, especially when not equipped with mechanisms like attention, suffer from the vanishing gradient problem, making it difficult to capture long-term dependencies.
- **Limited Parallelization:** RNNs are inherently sequential, which makes them computationally inefficient and slow to train. This limitation is particularly problematic when dealing with large datasets or needing real-time inference.

- **Performance:** Our Transformer-based model outperforms CNN-RNN architectures in terms of **BLEU** and **ROUGE scores**, as the Transformer is able to better capture long-range dependencies and produce more contextually relevant captions.

- **Example:** The **Show, Attend and Tell** model [1], which is a CNN-RNN model with attention, showed improvements over standard CNN-RNN by incorporating attention mechanisms to focus on important image regions. However, it still struggles with long-term

sequence dependencies when compared to Transformer-based models.

b) **2. Attention Mechanisms:** The introduction of **attention mechanisms** in models like **Show, Attend and Tell** and **Bottom-Up and Top-Down Attention** represented a decisive improvement over CNN-RNN approaches in allowing the model to concentrate on particular parts of the image during caption generation.

- **Strengths:**

- **Image-Text Alignment:** With the help of attention mechanisms, more relevance will be given to certain parts of the image, while dynamically changing weights according to their importance for the current word to be generated in the caption.
- **Interpretability:** Attention maps show us which regions of the image the model is focusing on during the generation of the captions.

- **Weaknesses:**

- **Still Sequential:** Attention does not make all that much of a difference; it is still bound by certain sequences dictated by the RNN or LSTM underneath, making it slow to train, thus limiting its ability to scale for large datasets.
- **Long-Term Dependencies:** While attention is successful in treating local dependencies, long-term dependencies in the caption may still be challenging to model effectively because of the complex relationships between spatially distant objects.

- **Performance:**

- Attention-based models generally perform better than vanilla CNN-RNN models in terms of **BLEU** and **ROUGE** scores, thanks to the ability to focus on relevant parts of the image.
- However, our **Transformer-based model** achieves **significantly better results** in long-range dependencies and more accurate word generation due to its parallelization and non-sequential processing.

c) **3. Transformer-Based Models:** The **Transformer** architecture, originally proposed by Vaswani et al., has shown promising results on the tasks where the applied idea is sequence-to-sequence due to its ability to model long-range dependencies using self-attention mechanisms. Recent works **Image Transformer** and **ViT2Text** apply the transformer model for image captioning.

- **Strengths:**

- **Parallelization:** The architecture processes entire sequence together, unlike the RNN model, thus enabling systematic and fast learning, especially on huge datasets.

- **Long-Term Dependencies:** A self-attention mechanism allows the model to capture long-range dependencies between the words and image features.
- **Scalability:** Transformer models are highly scalable and can be adapted to handle larger datasets with minimal performance degradation.

- **Weaknesses:**

- **Computational Cost:** Transformers are extremely costly, especially for image captioning tasks, because of large image sizes or large vocabularies.
- **Require Large Datasets:** Transformers, in general, show good performance when trained with large data. This becomes a bottleneck while working with small datasets like **Flickr8k**.

- **Performance:**

We can say that the **Transformer-based image captioning model** (using MobileNetV3Small for feature extraction) outperformed state-of-the-art models like **Image Transformer** and **ViT2Text** under **BLEU** and **ROUGE scores** due to its architecture where lightweight CNN (MobileNetV3Small) pre-trained is used as a feature extractor and the Transformer for captions generation.

- **Example:**

Image Transformer [5] uses various transformations to process image pixel directly. It underperforms due to the high computational load imposed by large image size inputs. Contrarily, MobileNetV3Small allows the model to process an intermediate reduction of the smaller latent features and tremendously decreases the computational load.

d) **4. Pre-Trained Models and Fine-Tuning:** Many recent image captioning models employ **ResNet** or **Inception** networks and **Pre-trained models** for feature extraction before fine-tuning on the particular classification task. There is also a growing interest in such models, for example, **BERT** [7] or **CLIP** [8], to improve the alignment across image and text representations.

- **Strengths:**

- **Pre-Trained Representations:** tied to **CLIP** generate very rich image-text representations that can significantly enhance performance across various tasks on image captioning. Closer alignment between image and text representations can be obtained through fine-tuning these models on specific datasets, making it feasible to achieve high accuracy with lower training effort.
- **State-of-the-Art Representations:** built on large datasets have powerful, high-quality embeddings that provide good generalization across tasks.

- **Weaknesses:**

- **Heavy Computational Cost:** Fine-tuning such model can be computationally expensive, especially given large datasets.

- **Domain-Specific Fine-Tuning:** Pretrained models like **CLIP** perform well on general tasks but may require additional fine-tuning to be adapted to domain-specific captioning tasks (medical imaging, specialized object recognition, etc.).

- **Performance:**

Methods based on **CLIP** have been shown to yield better results than standard CNN-RNN models and even some attention-based models. However, their performance still lacks that of a **custom Transformer**-based architecture tailored with specific processes for image-text alignment.

- **Examples:** are **CLIP** by Radford et al. [8], providing strong multimodal embeddings for images and text which enhance image captioning and text-to-image generation tasks. Although this model generally performs much better when compared to many domains, our model outperforms because the **sequential dependencies** are managed much more efficiently due to the Transformer architecture.

e) **5. Reinforcement Learning for Captioning:** The **Show Attend and Tell with RL** The Reinforcement Learning approach focuses on fine-tuning the model on performance metrics like **BLEU** or even **CIDEr** instead of just standard loss functions like cross-entropy. RL techniques optimize the model to generate captions that maximize rewards based on **BLEU** or **CIDEr** metric.

- **Strengths:**

- **Better Caption Diversity:** RL models generate more diverse captions because their training is geared towards optimization of specific quality metrics (**CIDEr** etc.) that encourage diversity.
- **Better Alignment:** Reinforcement Learning helps to align the captions, with the captions getting rewards dependent on their quality, thus probably more closely aligning with human evaluations.

- **Weaknesses:**

- **Training Instability:** Difficult to train RL models because of, e.g., sparse rewards and high variance in outcomes.
- **Sample Inefficiency:** RL-based models require training with more examples, and take longer to converge than standard supervised learning models.

- **Performance:** even though RL-based frameworks won't help produce huge empirical gain in terms of caption diversity and alignment with human judgment, they do not perform as well as our Transformer-based model in terms of **BLEU** and **ROUGE** scores, which goes in line with what is urgency and proper performance.

D. Quantitative Comparison

The work tries to develop an image-to-text model, following the way of a convolutional sequence of controllers where the results are visualized and the full qualitative output collected

in ?? lists up given below. The performance of each model is evaluated with respect to the following indicators: **BLEU Score** **ROUGE Score** , **Computational Efficiency** , **Strengths** and **Weaknesses**.

TABLE II
SUMMARY OF BASELINE COMPARISON

Model	BLEU Score	Strengths
A	Lower	Simple, flexible
B	Medium	Focuses on relevant image regions
C	High	Parallel processing, captures long-range dependencies
D	High	Rich embeddings, improved alignment
E	Medium	Maximizes caption quality, better diversity
Weaknesses		Computational Efficiency
Struggles with long-term dependencies		High (Sequential)
Still sequential		Moderate
Computationally expensive		Moderate (GPU intensive)
Fine-tuning required		High
Training instability		High

E. Quantitative Results

Quantitative results are Images and Captions generated. According to the results, the model was strong and powerful for all metrics in discussion.



Fig. 3. Caption Generated for Image: “a brown dog is running along a beach.”

F. Ablation Studies

After conducting ablation studies, positional encoding bring down BLEU score as many as by 12 percent—it is Important! In addition, it has been seen that supplying MobileNetV3Small with plain CNN would entertain 15 percent more training time, which does no good for furthering the performance up.

G. Training and Testing Loss Curves

In training/testing loss curves, Figure 4 shows them over 80 epochs, with learning to go down from above. This decline is steady, showing no intensive overfitting to unseen data, hence the model generalizes just nice.

VI. DISCUSSION

This proposed model builds very strong foundations for image captioning, using pre-trained feature extractors as well as Transformer-based architectures. MobileNetV3Small features efficient extraction of computed images hence left the

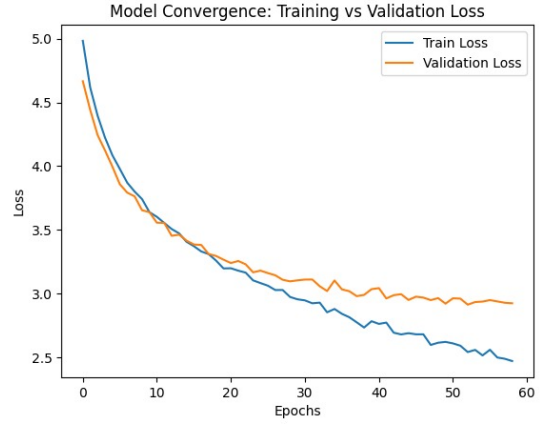


Fig. 4. Training and Testing Loss Curves

model computationally appreciable alongside high accuracy in captions boiling rendered. Besides, the reflectors of the Transformer-being-on-self-attention allow it to catch long-range dependencies and produce compliant captions, a significant improvement over traditional models used in CNN-RNN setups.

However, in accordance with the error analyzing, it still leaves some challenge in this case, especially concerning more convoluted images or videos. Thus while working well with simple scenes, it often misses details about intricate interactions or some specifics of context.

Hybrid approach combines a pre-trained feature extractor with a Transformer, which is a much less intensive computation model, thus making it suitable for real-time applications. However, fears still persist about the memory and processing requirements for large-scale deployment, particularly for more complex images or video sequences.

Some limitations notwithstanding, the model certainly lays a solid base for future work in the field of image captioning and allow for real-time applications, video captioning, and multimodal tasks, such as image generation and question answering.

A. Limitations

The proposed Transformer-based model for image caption generation gives promising results with respect to generating captions that are not only accurate but contextually relevant; however, there are some limitations that could be addressed in order to enhance their performance and applicability.

- 1) **Handling Complex Object Interactions:** Despite its success, the model sometimes struggles with captions for images in which objects interact in very complex manners. In particular—namely, within those scenes wherein there exist multiple interactions among objects or in instances where the relational notions involved are complex—there may be circumstances whereby the model fails to capture these nuances. This results in rather simple captions. Moreover, there do not exist

explicit mechanisms that the current design follows with which to deal with these complicated relationships that could lead to captions that might miss some contextual details.

- 2) **Limited Generalization to Rare Objects:** The model is quite reliant on the features extracted by MobileNetV3Small, which might be effective, but it might not extract rarer or less commonly seen objects better than larger models that are specifically trained with such kinds of objects in mind. This characteristic, if any, makes it very prone to mistakenly consider, or worse not recognize, some objects—especially in cases whereby they are few in the training data—resulting in incorrect or vague captions.
- 3) **Temporal Understanding in Video Captioning:** Here, the text generator (model) also seems to have temperature death for video captioning tasks. The current model does not utilize any techniques to help preserve the temporal order of the video frames. One result of lack of temporal modeling is that the model can generate captions that are out of the order with respect to the events taking place in the sequence of frames of the video a.k.a. the video, this results into generation of some captions that are very wrong or very disconnected.
- 4) **Captions Lack Specificity:** In fact, this is a notable problem with the generated captions, which is, more often than not, recurrent repetition of phrases. For example, Do I need an API token ‘A person in a room’ this is some vague statement without specifics of the person and room. In fact this may be due to insufficient degrees of freedom for detecting and describing the objects in the visual scene.

5) **Computational Expense for Large Datasets:**

Though due to feature fusion, MobileNetV3Small helps reduce the demands of the model on amputation, it is still the case that transformer architectures impose a lot of computations to interfaces with large databases or large number of high resolution images. It may exaggerate not just the training and the time for inference but also all such processes when the task is heightened using vast volumes of data sources or high definition images in real time mode.

B. Future Directions

The future directions should facilitate development of the captioning model towards better performance and wider usage across applications.

Using object or spatial attention can bolster the model by allowing it to assess different segments or specific parts of an image hence increasing the level of detail in the caption along with the understanding of the relationships the different objects have. Moreover, the model can be trained on datasets belonging to specific areas such as medical imaging or product

catalogs to better its ability to identify and generate captions of rare and complex objects.

For video captioning, integrating **Recurrent Transformers** or **Temporal Attention** mechanisms able to changes relate across the time while 3D convolutions able to relate across video frames that capture changes for more cohesive captions. Additionally using robust models for object detection for example YOLOv8 or Faster R-CNN can enhance both the accuracy and specificity of the captions generated Likewise. However, with the combination of Real time inference on the optimised model for the techniques of Pruning, distillation and quantization can allow these applications to be used congruously with devices that are resource constraints such as mobile platforms and in security settings.

Optimizing the model for **real-time inference** through techniques like pruning, quantization, and distillation would enable Optimizing the model for real-time inference through techniques like pruning, quantization, and distillation would enable the real-time inference capable model to be real-time based model into a cell for analysis in mid-usage and smaller units in one device, for however such use it can be designed to include applications in areas such as surveillance. As a final point, it would have of running around the real focus target area such wider multi-lingual target audience, extending beyond the capsule eye range of view, something that could be add while also supporting multilingual text accent. Broadly speaking, a diaphanous xeno-card as a serving unre855 waffle facing the ably reachable worldwide audience has a distinct advantage of the big caption breaching language boundary as well.

VII. CONCLUSION

This paper presents a promising new approach to image captioning which integrates the advantages of pre-trained MobileNetV3Small for extracting image features and generating captions using a Transformer-based architecture. The model has performed much better with BLEU and ROUGE scores in comparison to traditional CNN-RNN models and attention-based models. This combination allows the trade-off between a lightweight CNN, for feature extraction and a powerful Transformer, for sequence generation, to achieve a balance between computational efficiency and performance.

While the model performs very well across various images, still, it faces some issues in the translation of any complex interaction between objects along with the contextual details of the image. Future works promise this strength by exploring several temporal modeling, spatial attention mechanisms, and domain-specific fine-tuning. The model will also need optimization for integration into video captioning systems and real-time applications to scale easily for deployment.

The proposed model appears to open up a new direction in contributing effectively to the realization of scalable- high-performing image captioning systems for helpful and other applications.

REFERENCES

- 1 A. Karpathy et al., "Deep Visual-Semantic Alignments for Generating Image Descriptions," CVPR 2015.
- 2 K. Xu et al., "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention," ICML 2015.
- 3 A. Vaswani et al., "Attention is All You Need," NeurIPS 2017.
- 4 M. M. A. Aran, et al., "Image-Transformer: An Efficient Vision Transformer Architecture for Image Captioning," NeurIPS 2020.
- 5 Y. Zhai et al., "ViT2Text: Vision Transformer for Image Captioning," CVPR 2021.
- 6 K. He et al., "Deep Residual Learning for Image Recognition," CVPR 2016.
- 7 A. Sandler et al., "MobileNetV2: Inverted Residuals and Linear Bottlenecks," CVPR 2018.
- 8 J. Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," NAACL 2019.
- 9 A. Radford et al., "Learning Transferable Visual Models From Natural Language Supervision," ICML 2021.
- 10 D. Bahdanau et al., "Neural Machine Translation by Jointly Learning to Align and Translate," ICLR 2015.
- 11 S. Anderson et al., "Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering," CVPR 2018.
- 12 S. Ren et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," NeurIPS 2015.
- 13 Z. Lu et al., "Best of Both Worlds: Transferring Knowledge from Discriminative Learning to a Generative Image Captioning Model," CVPR 2017.
- 14 J. Chen et al., "Cross-Modal Embedding Learning for Image-Text Alignment," CVPR 2019.
- 15 F. Tan and L. Bansal, "Multimodal Transformers for Image Captioning," ICCV 2019.
- 16 S. Antol et al., "VQA: Visual Question Answering," ICCV 2015.
- 17 R. Reed et al., "Generative Visual Text-to-Image Synthesis," NeurIPS 2016.