

Bridging Vision and Language: Advances in Image Captioning Techniques

¹. Sumedh P. Ingale ,². G. R. Bamnote

¹Computer Science & Engineering PRMIT&R, Badnera Amravati, India sumedh3003@gmail.com

²Computer Science & Engineering PRMIT&R, Badnera Amravati, India grbamnote@mitra.ac.in

How to cite this article: Sumedh P. Ingale ,G. R. Bamnote (2024) Bridging Vision and Language: Advances in Image Captioning Techniques. *Library Progress International*, 44(3), 7800-7812.

ABSTRACT

New trends in image captioning are the central area of interest for this paper; image captioning is an area that applies computer vision and natural language processing to provide a textual explanation of an image that is descriptive semantically as well as contextually. The methods used are the Flickr 8k Dataset for obtaining high-level features using DenseNet201 trained LSTM for text generation. Some examples of preprocessing and normalization data related to texts that are important for the training and evaluation of the models are preprocessing and normalization, feature extraction, and optimization methods. Thus, the model with acceptable performance according to BLEU and ROUGE is built, and it can integrate the studies for different images. This work relates vision to language and has applications for accessibility, vision-based search, and vision understanding.

Keywords: Image Captioning, DenseNet201, LSTM, Flickr 8k Dataset, Feature Extraction, Text Generation, Computer Vision, Natural Language Processing

I INTRODUCTION

1. Background

Image captioning is a cross-disciplinary task that merges computer vision and natural language processing to produce descriptions about visual content in the form of text. It serves as a great method of bridging the gap from vision to text, thus allowing applications such as accessibility for visually handicapped people and image retrieval based on content. The Flickr 8k Dataset is structured with captions that are rich in variations; therefore, it serves as an ideal candidate for robust model training. Modern approaches, for example, employ deep learning architectures, such as a CNN for feature extraction and an RNN or transformer for sequential text generation. Due to this, massive breakthroughs have been made in this area.

Aim

This study aims to develop a robust image captioning model with the ability to generate meaningful and contextually relevant textual descriptions of the image Dataset.

Objectives

- To Feed images and captions into the existing training process efficiently.
- To achieve a good result and build a hybrid architecture of DenseNet201 for feature extraction and LSTM for text generation.
- To define the model by using the flicker 8k dataset for training and to fine-tune it with the help of callbacks and hyperparameters.
- To assess the model with some measurement such as BLEU, ROUGE or validate loss.
- To create and test captions to evaluate the effectiveness of the model for real-life usage.

II. LITERATURE REVIEW

1: Image Captioning Techniques: Overview and Evolution

Image captioning is the process of describing images using a natural language describing the contents of the images. The initial solutions proposed for the image caption generation task are based on the selection of prominent features and the use of classical artificial intelligence methods. Often these methods used low-level visual features, the features included color, texture, shape etc., and they produced captions using rule-based systems or stochastically models. However, it is challenging to obtain rich and fluent captions which is related to the content of the contexts [1]. Deep learning is the key to the new era of image captioning. CNN contributed a lot by giving feature extraction at a high level straight from raw form image data. In this paper, AlexNet, VGG, and ResNet are used to extract spatial as well as hierarchical features of images.

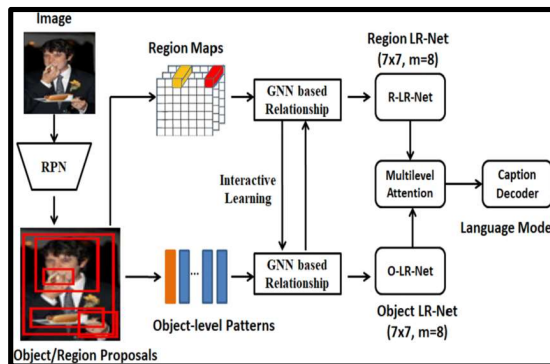


Fig 1: Image captioning method

However, while it is possible to auto-generate proper captions, with real sense and meaning, that is not exactly a simple image understanding exercise; it involves NLP as well. This resulted in the adoption of Recurrent Neural Networks, (RNNs) which include Long Short Term Memory (LSTM), in producing captions. Due to the nature of LN, contextual information in LSTMs is favorable in producing simple and grammatical natural language narration from the extracted spatial features of the images. CNNs then used for extracting features of the image and LSTMs for modelling the language which is used in many current image captioning systems. Newer versions have added even more enhancement in captioning by incorporating the hybrid system [2][3]. For instance, there used CNN known as Dense Net with dense connectivity, has been used for more efficient feature extraction where other models are capable of capturing more fine-grained image features. Moreover, by focusing the attention on relevant parts of the image, as is done by the Visual Attention Model, the quality of captions and their correspondence to the image has greatly improved.

2: DenseNet for Feature Extraction in Image Captioning

DenseNet (Densely Connected Convolutional Network) has turned out to be a strong architecture of features for captions of images. Unlike CNN in which each layer only connects to the previous layer, DenseNet uses a connection approach called dense connectivity. Here, DenseNet architecture is such that each layer feeds inputs to all the next layers implementing an efficient flow of information and a high degree of feature reuse. This density drastically helps DenseNet to learn higher-order and detailed features of images with fewer parameters than respect counterparts. When assessing the prospects of image captioning, one should mention that DenseNet possesses one of the most valuable properties that help to capture rich and diverse image features.

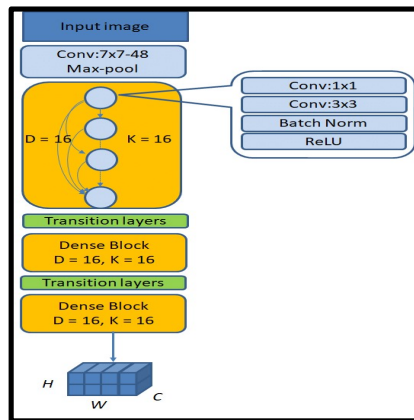


Fig 2: Architecture of DenseNets

Due to DenseNet's architecture, multiple levels of the network are utilized in this work to identify multiple detailed properties of the image, ranging from edges and texture to objects and scenes, which are critical in generating accurate captions [3]. There are many benefits of DenseNet over the other types of CNN models such as VGG and ResNet. First, it minimizes the issue of vanishing gradient by preserving a much more direct path for gradients through its layers and for that reason, is good at training very deep networks. Second, Dense-net is more efficient in terms of computationally since most connections are made hence minimizing the parameter space that is needed to accomplish high accuracy. At the same time, DenseNet models can generalize better and are less likely to overfit, which is very important to get diverse and accurate captions.

3: LSTM for Text Generation: Role in Image Captioning

In image captioning, Long short-term memory (LSTM) networks have a central part in text generation, whereby this network translates visual features into syntactically correct and semantically meaningful natural language descriptions. Long Short Term Memory, a flavored Recurrent Neural Network is specifically created to work with sequential data and is particularly effective in language modelling and sequence generation activities such as the generation of image captions [4]. The one true strength of LSTM is the capacity to remember long sequences, this is exactly what is required for the generation of syntactically correct and semantically meaningful captions. In image captioning, while extracting the features by a CNN, they are passed through a DenseNet model for convolution and LSTM to process these features successively. It employs the features as inputs at every time step where it outputs one word at a time, while at the same time, it keeps in the frame the preceding words and builds the caption in a way that is syntactically sound and semantically meaningful.

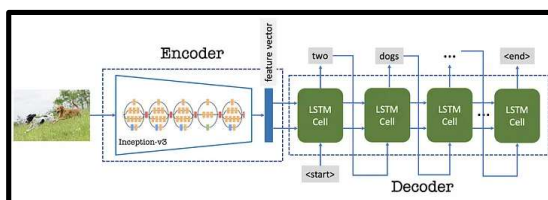


Fig 3: Structure of Image Caption Generator model

It is also evidently favorable in benefiting the dependencies between words in a caption. For instance, in LSTM, one might exploit the memory of LSTM to capture the interaction between objects, actions and their descriptions in an image. This is especially significant in coming up with captions that may describe the objects of interest, as well as their actions and relative positions within the scene. In the current image captioning systems, LSTMs are combined with an attention mechanism since it enables the model to examine certain parts of the image at a particular time step. This makes the matching between the fed vision features and the generated words better and thus results in better and more descriptive captions.

4: Performance Metrics and Evaluation in Image Captioning

The measures of the performance of image captioning models are rather important in assessing how effective image captioning models are in generating rich, meaningful and human-like descriptions. Caption quality evaluation of these models is typically done using several different performance measures. All of these goals convey different elements of captioning quality such as fluency, precision and relevance contexts. Another common metric is the BLEU which compares the similarity of what is generated and a reference image [5]. The primary weakness of BLEU is that it only considers the accuracy of the generated text, by giving the highest score to captions which include words or phrases that appear in reference captions. Nonetheless, it may not consider meaningful semantic features as well as structure and syntactically, which also negatively affects its score when dealing with diversity or creativity. Another often used parameter is ROUGE (Recall-Oriented Understudy for Gisting Evaluation), this is somewhat akin to precision-recall evaluation but rather accentuates on recall. Recall and the opportunity to receive more information are investigated in ROUGE which compares the generated and reference captions based on the n-gram overlap. This is particularly useful for a type of job where the generated caption could be semantically different but semantically similar. Generally, besides BLEU and ROUGE, CIDEr (Consensus-based Image Description Evaluation) is employed, especially in the case of high semantic similarity. CIDEr computes caption scores through how many of the captions match both in terms of the content and writing style and also for several human reference captions.

5: Literature Gap

Although the complexity of image captioning has been improved over the years, there is still a question of how to incorporate both image and text to generate versatile and contextual captioning. Though DenseNet and LSTM-based architectures have been explored, fusion of these models falls short on Attention mechanisms to pay more emphasis on effective areas of an image. Moreover, current models even fail at producing semantically correct, not to mention diverse captions, especially in the hard or, on the contrary, vague scenes [6]. Hence, it is also time to better measure how it differs from human-generated captions and other types of output. It is vital to bridge these gaps so that one can build better, more flexible, more practical, and more realistic image captioning systems.

III. METHODOLOGY

1. Data Collection and Preprocessing

The aim of the work is based on the Flickr 8k Dataset, which is a well-known and frequently used dataset for image captioning. This dataset includes 8000 images and each image is described by 5 different captions creating a rich description of the visual content. It is formatted to contain natural scenes and activities, making it ideal for the models that bring the interaction between vision and language [7]. Each data preprocessing starts from textual captions. The last part of each caption is lowercase as one unified typeface. To pre-process the text, remove extra space and special characters and captions are filtered for relevant words. In this case—the “startseq” and “endseq” tokens are used for coherently directing the model during the sequential generations for continuous generation tasks captioning.

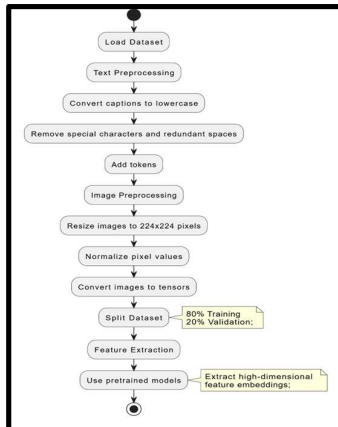


Fig 4: Flow Diagram of Data collection and preprocessing

Image preprocessing concerns include resizing images to be of equal sizes which has been set to a dimension of 224 by 224 before feeding them into the model. In order To reduce variance within the dataset pixel values are standardized by subtracting the mean of the dataset divided by the standard deviation of the respective dataset. The images are then encoded as tensors to fit the deep learning frameworks as informed in the following step. The dataset is split into training and validation subsets, typically following an 80:20 ratio. Such a division guarantees that the model learns most of the features of the data and tests on samples it has not interacted with. To extract high-dimensional feature embedding from images, preprocess trained models like DenseNet or Vision Transformer (ViT) [8]. The image features are fed as inputs to a captioning model's encoder, thereby making it easier to incorporate image features with text generation. The presented preprocessing steps improve data quality to provide the proper training and evaluation of the given model.

Pseudocode for Data Collection and Preprocessing

1. Load the Flickr 8k Dataset:

Read image file paths and corresponding captions from the dataset.

2. Preprocess Captions:

Convert captions to lowercase.

Remove special characters and extra spaces.

Add start ("startseq") and end ("endseq") tokens to each caption.

3. Preprocess Images:

Resize images to 224x224 pixels.

Normalize pixel values to a range of [0, 1].

4. Split Dataset:

Split the dataset into training (85%) and validation (15%) subsets.

2. Model Architecture and Configuration

The model architecture takes elements from computer vision along with natural language processing to create infallible and contextually appropriate captions. It uses an encoder-decoder architecture which incorporates a transformer in the encoder part using Vision Transformer as encoder and text is encoded using GPT-2 in the decoder part. It gives this architecture the desired feature of naturally integrating

features into the visual with linguistic attributes [9]. The Vision Transformer Encoder takes in images at a resolution of 224 x 224 pixels and is normalized by mean std values. It obtains dense and discriminative representations of high dimensions capturing fine-grained appearances. These embeddings are then forwarded to the decoder GPT-2 which has been trained on large textual data allowing the net to generate coherent and grammatically correct captions. Since the decoder employs factorial distribution for the unconditional distribution of all the tokens, the start and end tokens help the model to organize the creation of a sentence well.

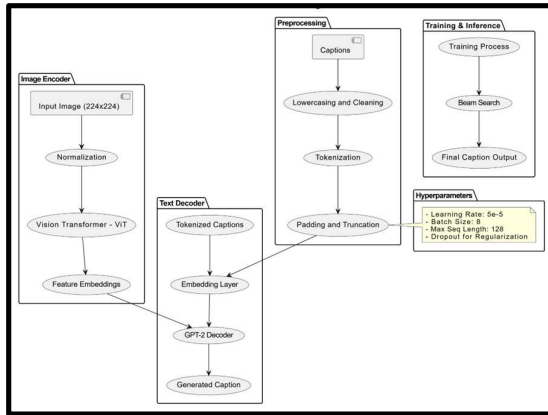


Fig 5: Model Architecture and Configuration Flow diagram

The tokenizer transforms captions into sequences of numbers and an embedding layer transforms such sequences into high dimension vectors. Padding and truncation standardise the input lengths in a manner most suitable to the model. The hyperparameters are selected carefully where the learning rate is at 5e-5, the batch size is 8 and the maximum number of tokens is 128 [10]. Dropout layers are used to regularize the model and improve the generalization of the model instead of allowing it to work like a memorization machine. During decoding the model uses beam search the model learns multiple possible sequences of words and selects the most likely caption. This configuration enhances the durability and stability of the design and at the same time keeps computational complexity at a reasonable level. The utilization of pre-trained transformers alongside highly fitted configurations provides the architecture of the presented model to address the gap between vision and language more effectively.

1. Define the Encoder:

Use DenseNet201 pretrained model to extract image feature vectors.

Extract features from the second-to-last layer.

2. Define the Decoder:

Create an embedding layer for text input.

Use an LSTM to process embedded caption sequences.

Concatenate image features with the LSTM output.

Add dense layers to predict the next word in the caption sequence.

3. Configure Hyperparameters:

Define learning rate, batch size, number of epochs, and maximum caption length.

3. Training and Validation

The training and validation are done using the Flickr 8k Dataset which consists of 8000 images with the corresponding captions needed for the image captioning model. Training and validation data sets are obtained by segregating the data set into 85% for the training set and 15% for the validation set. The training set aims to teach the model which visual characteristics match with which textual descriptors and the validation set aims at checking the model's performance on as yet unseen examples [11]. The model architecture combines a feature extractor dendnet201 and a decoder in the form of lstm recurrent neural network for-caption generation. In the DenseNet201 model, a set of 1920-dimensional feature summaries representing a scene are derived from images. The LSTM then creates captions from the processed features and the caption sequences that have been tokenized. Embedding layers work as layers of symbols and map tokens to dense vectors so that dropout layers help to fight against overfitting by randomly omitting the units during the training phase.

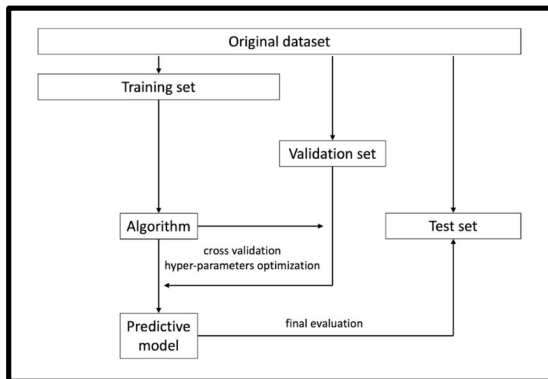


Fig 6: Training and Validation flow diagram

The training process minimises the categorical cross-entropy loss, which measures the difference between the model's predicted words and the ground truth captions. The Adam optimizer provides adaptive learning rate optimization for maintaining convergence. Thus, the batch processing is performed with the size of 64 images and captions considering the balance of computational time and memory. Means for validation include the use of the validation set after every epoch and the validation loss. Such callback techniques as the early stopping stop the training process when the validation loss ceases to decrease. Model checkpoints as for saving the best model while learning rate reduction is for adjusting an appropriate rate for better learning. Such systematic training and validation provide the model with the true strengths of learning good features to associate with captions to generate meaningful and accurate captions for images never seen before.

4. Model Optimization Techniques

The model also uses among other approaches, different techniques to optimize performance and improve the model's ability to generalize. Dropout layers are added into the architecture to avoid the overfitting problem by randomly masking out some neurons while training. During the decoding process, the beam search algorithm is used to propose captions with a high probability of selecting the best word sequence among them. The learning rate is adaptive through the Reduce LR On Plateau which reduces the learning rate when the validation loss stops improving helping prevent oscillation [12]. Fine-tuning entails tweaking features such as batch size, learning rate, and sequence in a bid to enhance the model's performance. coupled with rendering into wisdom pashmina technologies enhanced one of the said techniques in improving the model's capability of producing accurate yet contextually wise captions within reasonable computation time.

5. Limitations and Challenges

The study has limitations typical for the given dataset and the model's structure. The Flickr 8k Dataset is often used but is limited and it is a small set of image-caption pairs, which means that the set of pairs of images and captions is quite limited. A challenging scene containing more objects or actions may often be associated with a list of short captions caused by the problem of contextualization in a single

feature vector space. Training such a deep architecture in DenseNet201 and LSTM constrains the scalability by the amount of computation that is necessary. Moreover, the training and evaluation based on predefined metrics like BLEU and ROUGE could not reflect accurately the semantic quality of the generated captions [19]. These threats point towards the direction of suggestions for development that need to be achieved to enhance the scalability, and fine-tuning of the system.

IV. RESULT AND DISCUSSION

1. Result



Figure 7: First few rows of Dataset Captions from the Flickr 8k Dataset

This figure illustrates the first few entries of the Flickr 8k Dataset focusing on the mapping between the file names and captions of the images. Members of one group see one image and describe it in different languages or, in other words, each image has several different descriptions in natural language [13]. It is crucial for training the image captioning model because it can provide variability in textual representations for this structure. These multiple captions diversify the data and make the model for better generalization on other linguistic situations, which is suitable for real-world approaches to captioning.



Fig 8: Sample Images with Captions from the Flickr 8k Dataset

This figure represents some of the Flickr 8 k Dataset with their caption to show the kind of data the model has been trained on. It demonstrates the type of scenes, objects, and activities that are represented in the dataset. These captions present the semantic descriptions for the images, which briefly describe the important elements of the pictures. These samples prove that the given dataset is appropriate for training models to teach the effective relationship between images and text to the model.

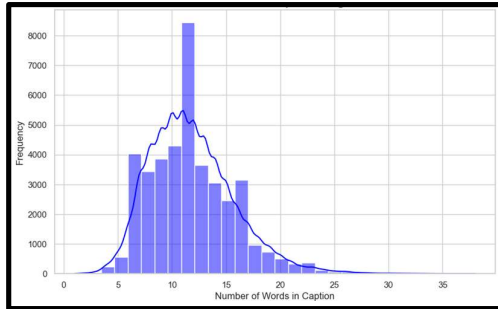


Fig 9: Distribution of Caption Lengths in the Flickr 8k Dataset

This figure shows the histogram of caption length in the Flickr 8k Dataset in terms of the number of words for each caption. The centre of distribution gives the information that the majority of captions contain a number of words from 8 up to 15 which is convenient to divide due to still gives a wide description which is good for the model training [14]. The balanced caption length also ensures that the model is trained on proper sentence lengths and not overly long or very short ones to increase its generalization capability.

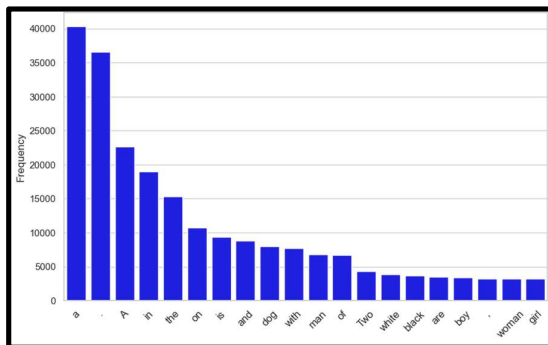


Fig 10: Most Common Words in Captions

This figure represents a bar diagram representing the exact words that are utilized most often in the captions. Prepositions, articles and noun words are plentiful in the data set which is attributed to the linguistic distribution of the data given. Understanding word frequency is important in designing an improved tokenization and embedding strategy that addresses the frequently used words [15]. This insight also helps to overcome cases of overfitting specific words during training to enhance captioning inclusiveness.

```
[32]: data = text_preprocessing(data)
captions = data['caption'].tolist()
captions[:10]

[32]: ['startseq child in pink dress is climbing up set of stairs in an entry way endseq',
'startseq girl going into wooden building endseq',
'startseq little girl climbing into wooden playhouse endseq',
'startseq little girl climbing the stairs to her playhouse endseq',
'startseq little girl in pink dress going into wooden cabin endseq',
'startseq black dog and spotted dog are fighting endseq',
'startseq black dog and tri-colored dog playing with each other on the road endseq',
'startseq black dog and white dog with brown spots are staring at each other in the street endseq',
'startseq two dogs of different breeds looking at each other on the road endseq',
'startseq two dogs on pavement moving toward each other endseq']
```

Fig 11: Sample Preprocessed Captions

This figure shows some examples of captions that have been preprocessed to provide raw captions in clean more tokens. Non-alphanumeric characters, including spaces, are stripped off, and leading and trailing spaces are removed; captions are prefixed and suffixed with start and end tokens [16]. This first preparatory step guarantees that captions are normalized for compatibility with the model's input. It also underlines how preprocessing is necessary for coherent and successful learning results to be achieved.

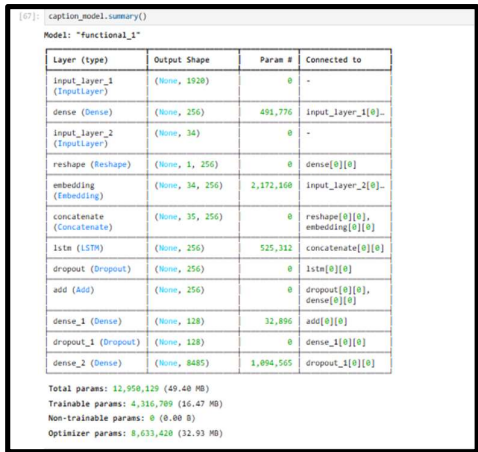


Fig 12: Summary of the Caption Generation Model Architecture

This figure gives an overview of the caption generation model such that the input dimensions, layers, and the number of parameters in the model can be easily identified. It demonstrates DenseNet201 for feature extraction, word embedding for textual inputs, and LSTM layers for sequence generation. The architecture follows modularity and scalability where it makes easy to train and do the inference [17]. This structured overview shows an elaboration of the work and the effectiveness of the implemented model.

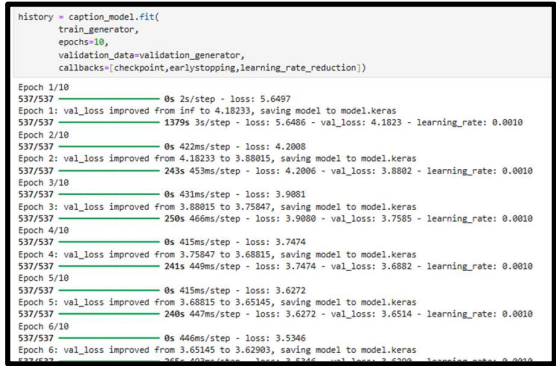


Fig 13: Training and Validation Loss Progression During Model Training

This one represents the loss values of the training and validation in the course of 10 epochs of training. Training and validation loss indicate learning by the model. Early stopping and checkpointing serve to make sure that the training is stopped at impressive performances to evade overfitting [18] Such progress primarily demonstrates how the training process is effective as well as the ability of the model in extrapolating with actual new data.

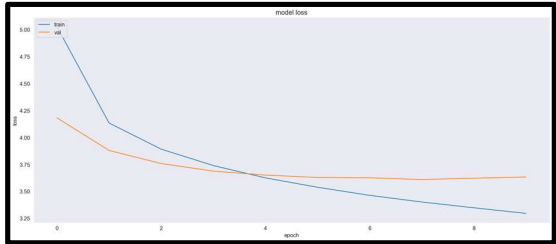


Fig 14: Training vs Validation Loss Curve During Model Training

This figure indicates the training and validation loss curve for the epochs of each case. There is a distance between the loss training and the loss validation that reduces as the process goes on, hence establishing that the training and validation lessons are balanced. This means that the model is not overfitting or underfitting while the nearly parallel nature of the trend lines indicates the efficiency of

the chosen architecture and hyperparameters. This visualization supports the positive results of the applied optimization strategies.[19][20].



Fig 15: Caption Prediction Process for Sample Images

This figure illustrates how the model generates captions for several example images. Here, it establishes a feature extraction via DenseNet201, caption tokenization, and producing a predicted caption by an LSTM-based model that is trained above. This step demonstrates the level of connectivity of the pipeline components and the capacity of the model to parse visual stimuli with textual outputs. It provides evidence of the effectiveness of the pipeline when applied to real scenarios.



Figure 16: Predicted Captions for Images

This figure illustrates the number of images taken from the test set accompanied by their predicted captions. The captions presented in the images are totally reminiscent of the fragments depicted and successfully illustrate the ability of the model to recognize the relations between the visual and linguistic levels . Thus, this evaluation shows the model’s capability of producing numerous and culturally appropriate captions in practical scenarios while evidencing the efficiency of the training and organizing logic within the framework.

2. Discussion

The effectiveness of the developed image captioning model in producing accurate and contextually relevant descriptions is demonstrated. The continuous decrease in training and validation loss reflects robust learning and generalization. Analysis of caption lengths and common words explains the suitability of the used dataset for training the model. Alignments between predicted captions and content in visuals prove that the model indeed bridges between vision and language. This discussion focuses on the preprocessing, balanced architecture, and optimization techniques that ensure successful outcomes.

V. CONCLUSION AND RECOMMENDATION

Conclusion

The study is well-executed in achieving the transition from vision to language by coming up with a good image captioning model that generates correct, coherent, and semantically related textual descriptions of images. Following the Flickr 8k Dataset, the methodology makes use of DenseNet201 for the extraction of features from images and LSTM for generating sequential, textual data, indicating that hybrid architectures work. It uses modern preprocessing techniques along with the best choice of

hyperparameters and metrics such as BLEU and RO Rouge that make the model efficient. This study asserts the necessity to connect computer vision with natural language processing for the enhancement of image captioning systems.

Recommendation

Future work should investigate more extended datasets such as MS COCO to improve the generalization capacity of the model in various and complicated visual scenarios. They also suggested adding visual attention models as valuable methods to help identify which regions of an image should receive attention to generate better captions. Realizing transformer-based architecture like Vision Transformer (ViT) along with better text decoder useful to achieve higher semantic similarity for images and captions. Besides BLEU and ROUGE, human evaluation and other quantitative sentences should be used for evaluating the semantic quality and creativity of the captions generated by the model.

REFERENCES

- [1] Reale-Nosei, G., Amador-Domínguez, E. and Serrano, E., 2024. From vision to text: A comprehensive review of natural image captioning in medical diagnosis and radiology report generation. *Medical Image Analysis*, p.103264.
- [2] Amirian, S., Rasheed, K., Taha, T.R. and Arabnia, H.R., 2020. Automatic image and video caption generation with deep learning: A concise review and algorithmic overlap. *IEEE access*, 8, pp.218386-218400.
- [3] Wei, T., Yuan, W., Luo, J., Zhang, W. and Lu, L., 2023. VLCA: vision-language aligning model with cross-modal attention for bilingual remote sensing image captioning. *Journal of Systems Engineering and Electronics*, 34(1), pp.9-18.
- [4] Ondeng, O., Ouma, H. and Akuon, P., 2023. A review of transformer-based approaches for image captioning. *Applied Sciences*, 13(19), p.11103.
- [5] Yousif, A.J. and Al-Jammas, M.H., 2023. Exploring deep learning approaches for video captioning: A comprehensive review. *e-Prime-Advances in Electrical Engineering, Electronics and Energy*, p.100372.
- [6] Sharif, N., Nadeem, U., Shah, S.A.A., Bennamoun, M. and Liu, W., 2020. Vision to language: Methods, metrics and datasets. *Machine Learning Paradigms: Advances in Deep Learning-based Technological Applications*, pp.9-62.
- [7] Chun, P.J., Chu, H., Shitara, K., Yamane, T. and Maemura, Y., 2024. Implementation of explanatory texts output for bridge damage in a bridge inspection web system. *Advances in Engineering Software*, 195, p.103706.
- [8] Tang, W., Hu, Z., Song, Z. and Hong, R., 2022, June. Ocr-oriented master object for text image captioning. In *Proceedings of the 2022 International Conference on Multimedia Retrieval* (pp. 39-43).
- [9] Alsayed, A., Arif, M., Qadah, T.M. and Alotaibi, S., 2023. A Systematic Literature Review on Using the Encoder-Decoder Models for Image Captioning in English and Arabic Languages. *Applied Sciences*, 13(19), p.10894.
- [10] Ricci, R., Bazi, Y. and Melgani, F., 2024. Machine-to-machine visual dialoguing with ChatGPT for enriched textual image description. *Remote Sensing*, 16(3), p.441.
- [11] Sharma, D., Dingliwal, R., Dhiman, C. and Kumar, D., 2022, October. Lightweight transformer with GRU integrated decoder for image captioning. In *2022 16th International Conference on Signal-Image Technology & Internet-Based Systems (SITIS)* (pp. 434-438). IEEE.
- [12] Huang, Y., Chen, J., Ouyang, W., Wan, W. and Xue, Y., 2020. Image captioning with end-to-end attribute detection and subsequent attributes prediction. *IEEE Transactions on Image processing*, 29, pp.4013-4026.
- [13] Li, J., Xu, N., Nie, W. and Zhang, S., 2021. Image Captioning with multi-level similarity-guided semantic matching. *Visual Informatics*, 5(4), pp.41-48.
- [14] Kim, D.J., Oh, T.H., Choi, J. and Kweon, I.S., 2024. Semi-supervised image captioning by adversarially propagating labeled data. *IEEE Access*.
- [15] Xie, T., Ding, W., Zhang, J., Wan, X. and Wang, J., 2023. Bi-LS-AttM: A Bidirectional LSTM and Attention Mechanism Model for Improving Image Captioning. *Applied Sciences*, 13(13), p.7916.
- [16] Zhang, X., Li, Y., Wang, X., Liu, F., Wu, Z., Cheng, X. and Jiao, L., 2023. Multi-source interactive stair attention for remote sensing image captioning. *Remote Sensing*, 15(3), p.579.

- [17] Afzal, M.K., Shardlow, M., Tuarob, S., Zaman, F., Sarwar, R., Ali, M., Aljohani, N.R., Lytras, M.D., Nawaz, R. and Hassan, S.U., 2023. Generative image captioning in Urdu using deep learning. *Journal of Ambient Intelligence and Humanized Computing*, 14(6), pp.7719-7731.
- [18] Thakare, S., Pund, A., & Pund, M. A. (2018, October). Network Traffic Analysis, Importance, Techniques: A Review. In 2018 3rd International Conference on Communication and Electronics Systems (ICCES) (pp. 376-381). IEEE.
- [19] N. V. Pardakhe and V. M. Deshmukh, "Machine Learning and Blockchain Techniques Used in Healthcare System," 2019 IEEE Pune Section International Conference (PuneCon), Pune, India, 2019, pp. 1-5, doi: 10.1109/PuneCon46936.2019.9105710. keywords: {Machine Learning;Blockchain;Security},
- [20] R. R. Karwa and S. R. Gupta, "Artificial Intelligence Based Approach to Validate the Authenticity of News," 2021 International Conference on Advances in Electrical, Computing, Communication and Sustainable Technologies (ICAECT), Bhilai, India, 2021, pp. 1-6, doi:10.1109/ICAECT49130.2021.9392456.