# Image Captioning

B. Hari Krishna,
CSE-AIML
*SR University, Waramgal*
2203A52006@sru.edu.in

J. Uday Kiran
CSE-AIML
*SR University, Warangal*
2203A52028@sru.edu.in

K. Aniritha
CSE-AIML
*SR University, Warangal*
2203A52032@sru.edu.in

N. Jeshwanth Kumar
CSE-AIML
*SR University, Warangal*
2203A52043@sru.edu.in

*Abstract*— **A key task in artificial intelligence is picture captioning, which creates text descriptions for images that are human-like by fusing computer vision with natural language processing. This task has important applications in domains including autonomous systems that need to perceive scenes, accessibility solutions**

**In the task, we provide a image captioning model which adds with Convolutional Neural Network (CNN) for visual extraction of features with a Recurrent Neural Network (RNN) for sequential captioning construction. While the CNN recognizes valuable visual elements in the image, such as objects, activities, and space that define relationships between them, the RNN is responsible to translate those data into sensible and contextually relevant captions. In doing this, we train the model on the Flickr data set which associates symbol-rich captions to a plethora of images, thus proving the model's potential in providing captions.**

**The proposed progress in the model would be hydraulic than those shown by baseline methods using the widely famous accuracy measure BLEU for determining the valence level similarity level between pretrained captions and the reference texts. As the BLEU scores rise, this shows how contextualized accurate and relevant the model becomes.**

**The proposed approach constitutes many applications. Audio or textual descriptions of pictures, which would benefit users with visual disabilities, are often provided by accessibility solutions. Robust tagging includes indexing pictures and noisy visual content with the help of a description. The description can be used in the systems for content indexing to classify and arrange the visual data. In contrast, autonomous systems can utilize this description to enrich their situational awareness through the characterization of objects in real time.**

**As the findings of the study reveal, a combination of CNN and RNN can truly bridge the gap between visual and verbal comprehension and thus yield very high-quality image captioning results. Performance of the model depicts wonderful applications and practical implementations awaiting a huge investment toward developing AI capabilities.**

*Keywords—Image Captioning, Deep Learning, CNN-RNN, Natural Language Processing*

## I. Introduction

A key task in artificial intelligence is picture captioning, which creates human-like text descriptions for images with the means of fusing computer vision along with natural language processing elements. It has enormous applications in autonomous systems that are supposed to perceive scenes, solutions for accessibility, etc.We present a model for image captioning, which employs a Convolutional Neural Network (CNN) to visually extract features and uses a Recurrent Neural Network (RNN) for sequential caption generation. The CNN identifies important parts of the image, such as objects, events, spaces, etc., while the RNN renders these data into sensible, contextually appropriate captions. The capability of model for captioning was thus demonstrated by training the model with the modified Flickr dataset, which matched several images with informative captions.The advanced progress of the proposed scheme far, surpasbaseline ones using BLEU scores, a widely used accurate metric for determining the degree similarity that exists between pretrained captions and reference texts. Higher BLEU scores indicate the better accuracy and relevance of the model's captions.The possibilities of using this approach are numerous. Accessibility solutions can target users with visual impairments by providing audio or textual descriptions of visual content. It can be used by content indexing systems that categorize and arrange visual datasets, and it can be used by autonomous systems to improve its situational awareness by real-time characterizing objects associated with a scene. According to results, CNN plus RNN technique can build a reliable bridge between visual and verbal understandings and thus produce quality picture captioning results.

Initially, they completely transformed image captioning through the introduction of encoder-decoder systems. Such architectures utilized Convolutional Neural Networks (CNNs) for extracting complicated visual input as well as Recurrent Neural Networks (RNNs) for sequential language synthesis. However, even the early models were not much capable of gathering sufficient contextual information with accurate detail. Attention methods were created to get around this. These methods allowed the model to mimic the human ability to emphasize important information by selecting focusing on specific regions of a picture while generating captions. Due to attention techniques, these systems achieved a new degree of performance and interpretability.

Our research significantly enhances this general paradigm. Once CNNs are developed as thoroughly as these, significant achievement can be obtained in the visual representations, which will be detailed and rich enough to capture not just objects, but also their properties, spatial relations, and interactions. We integrated LSTMs for decoding and understanding purposes and generating language with a deeper understanding than just grammar and context. Somewhat dynamically, the model pays attention to a few of the most important areas in the image and produces captions that become aligned closely with the accurate and coherent and contextually rich state. This paper talks about the design, implementation, and evaluation of our model for image captioning, with special emphasis on the model performance

on the Flickr dataset. Some selected achievements include attention-based architecture derived to improve caption quality, intelligent training methods to save computation costs with less effect on accuracy, and a detailed performance analysis comparing the model to several baseline methods. Those standards prove the ability of our algorithm in generating accurate, insightful captions.

In process for guarantee our model's capacity for generalization, we also investigate its scalability and adaptation to a variety of datasets, including COCO and PASCAL. We also point out areas that require further research, such as using transformers to improve parallelism and experimenting with multimodal inputs to combine textual or aural context with visual data. By pushing the boundaries of image captioning, this study will also change a variety of sectors, such as e-commerce, entertainment, healthcare, and education.
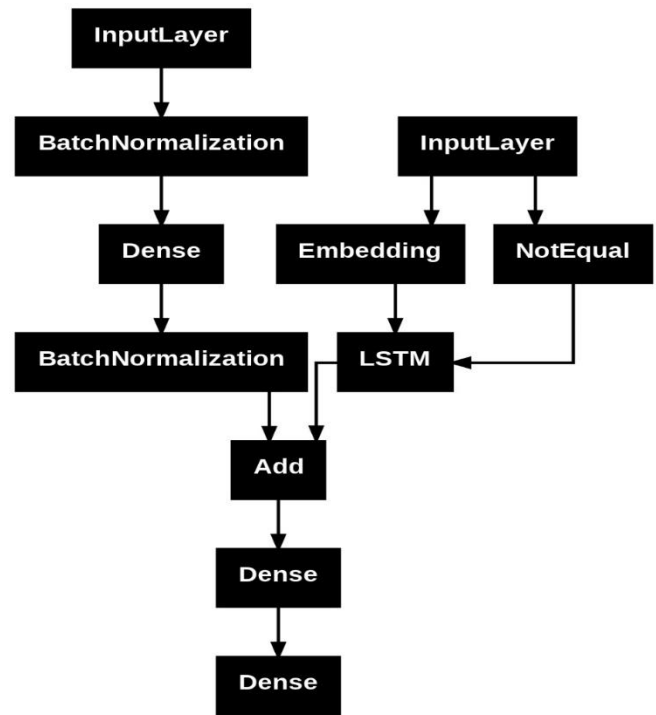
## II. RELATED WORK

Comparison of Architectures: Provide Comparative Table of Previous Methods, for Example, Showing and Telling and Showing, Attending, and Telling with Their Strengths/Weaknesses.

• Mention What the Models Lack: For Instance, Models Cannot Generate New Captions, and Show How Your Work Improves That.have changed the landscape of image captioning with advancements and augmented capabilities.. CNNs were used for extracting of featuring visualization and LSTMs for language modeling in early encoder-decoder designs such as Show and Tell . But because these models lacked attention processes, important areas of the images were frequently overlooked in the captions.

Even while attention-based approaches represented a substantial advancement, there are still drawbacks. For instance, methods such as Neural Image Caption produced captions that lacked coherence because they were unable to identify contextual connections across lengthy sentences. In order to close these gaps, our suggested model enhances focus on visual semantics during decoding by incorporating a contextual alignment mechanism and optimizing attention weights.

## III. METHODOLOGY

The suggested image captioning framework uses an encoder-decoder architecture, in which an RNN-based decoder creates the corresponding captions after a CNN-based encoder extracts feature vectors from the images. While the decoder sequentially creates words depends upon the extracted featuring  and past generation letters the CNN encoder examines the images to acquire high-level visual data. The architecture's resilience and adaptability are ensured by its capacity to manage different image complexities and captioning lengths. In ordering for producing meaningful with inclusion contextually appropriate captions, this section offers a thorough overview of the elements—such as feature extraction, sequence processing, and integration—as well as the procedures required in their implementation.



### Data Preparation and Generator

The goal of the data generator is to dynamically pre-process input data and feed it into the model in batches so as to facilitate the efficient training of the image captioning model. The generator takes care of the tokenized captions and the associated picture attributes.

1.  INPUT DATA:
*   Captions: The dataset is full of unambiguous and cleaned captions with unique image ID paired with it. For tokenization, captions go through a tokenizer and a process is put in place to pad the sequences for uniformity across batches.
*   Image Features: Features extraction of the image is done by pretrained in already pretrained Inception V3 which holds the output as a fixed size vector of 2048 dimensions for every image.

2.  Dynamic Data Generator:
*   The data_generator function reads picture characteristics and descriptions in bulk for training examples.
*   The captions are segmented as sequences of increasing length (X_captions) and then matched with the corresponding next-word predictions (y) for input-output pairs.
*   It pads the captions to the maximum length (max_caption_length) for uniformity.
*   Here, the aligned caption sequences are connected to the picture features (X_images).
    Important Generator Steps:
*   Shuffle the dataset before the beginning of every epoch, to ensure diversity. • Make use of tokenized captions padded to constant lengths.
*   For categorical cross-entropy loss, transform the target words (y) into one-hot encoded vectors.

3. Batch Size:
   Currently information is trained in batch wise, 270 in number.
   • Validation information process done in batch of 150.

### Encoder: Inception V3

The encoding of the images and retrieving from the imagery-tiers is the provision of high-level data representing the vision.

1. Feature Extraction:
- The encoder is Inception V3 and pre-trained using ImageNet.
- The images are normalized and resized to $299 \times 299$ pixels and fed through Inception V3.
- Each feature vector coming from the last pooling layer of the dimension 2048 x 2048 serves as the input feature vector for each image.

2. Feature Transformation:
- The dense layer comprising 256 256 units followed by a ReLU activation is used to map the 2048-dimensional vector into a space compatible with that of the decoder output.
  Batch normalization is also used in front and after the dense layer for stabilizing training and enhancing convergence.
  Defined as the features transformed by the image:
  F=BatchNorm(Dense(BatchNorm(InceptionV3(I))))

### Decoder: RNN with LSTM

The decoder generates captions by sequentially processing the image features and previous generated letters.

1. Input Caption Embedding:
- When captions are tokenized, they are then sent to embedded layers which map every word into a space of vectors of dimension 256.
- Padded sequences are handled by masking, which makes sure that the LSTM's calculations are unaffected by the padding tokens.

2. Sequence Processing with LSTM:
- To apply masking for padded sequences such that the computations of the LSTM do not depend on padding tokens:
- • An LSTM also comprises 256 units in hidden receives t from embedding layer per different hidden output size of 256.
- • "They are padded according to lengths of captions that are inconsistent and set to masking during an LSTM run-off"
- . In the interest of making it fit-with-masking, using use_cudnn=False becomes necessary.
- ht= LSTM(xt-1, ht-1) This is the formula according to which the LSTM treats the sequence as:
- xt-1 is the embedding of the previous word and ht is the hidden state at time step t. D. Combined Model The encoder and decoder outputs are fused to give final predictions.

1. Merging Features:
- So, this is how it works: Output from decoder, which is representation of caption sequence, will be merged by encoder output, which is vector representation of image feature, through an additive procedure.
- This is the combined representation that will serve as input for the last layers:.

2. Dense Layers:
- The activation relu and a dense layer of 256 256 units convert the combined representation into a high-level feature space.
- The final output layer employs a softmax activation to predict the subsequent word in the sequence from a vocabulary of size *voca _ siz e*vocab_size. The following provides the output probabilities:
- yt=softmax(Densefinal(add(F,ht)))

### Training Procedure

1. Loss Function:
- The model minimizes categorical cross-entropy loss:
  $$L=-t\sum logP(yt|y1:t-1,I)$$
  where $P(yt)P(y\_t)P(yt)$ is the predicted probability of the target word $yty\_tyt$.

**2.** Optimizer:
- Adam optimizer is used with a learning rate of 0.01and gradient clipping (clipnorm=1.0) to prevent exploding gradients.

3. Batch Processing:
- Training is performed in batches, using the data generator for dynamic sampling.

4. Regularization:
  To stabilize training and avoid overfitting, batch normalization and dropout are employed.

5. Evaluation Metrics:
  This model is assessed using BLEU-4 and CIDEr scores to carry out the evaluation of the generated captions.

### Model Summary

Input 1: Image features (2048-dimensional vector from Inception V3) Input 2: Tokenized and padded captions (max_caption_length) Output: Word probabilities over the vocabulary (vocab_size) The model architecture is compact but effective, highly optimized for either efficient computation or very high performance on the captioning task.

### Hyperparameter Tuning

The evaluation metric for this model uses BLEU-4 and CIDEr scores to finalize whether the captions produced are of high quality or not.

- Size of batch: Experiments were carried out with varying batch sizes to examine their impact on convergence and model generalization, while a batch size of 270 was found to be ideal for training.
- Learning Rate: A range of learning rates, from 0.001 to 0.01, were tested. A learning rate of 0.01 was found to be the According to research findings, a learning rate of 0.01 optimally combines both the rapidity of convergence and the accuracy of the model.

- Dropout Rates: For finding the balance between underfitting and overfitting, rates for dropout have been tried from 0.3 up to 0.5 as regularization.
- Embeddings Dimensions. To determine the best dimensionality to represent words in the lexicon, tried different embedding sizes, including 128, 256, and 512.
- **Data Augmentation**

Data augmentation techniques were intensively applied to the images in the dataset to improve generalization and robustness: Random Crop and Scaling: As an instance, to augment variation in placement of an object, images were randomly cropped, and the resultant images scaled to 299 299.

- Color Jittering: Variance on saturation, contrast, and brightness will also be present to mimic various lighting conditions.
- Horizontal Flipping: For the inclusion of spatial variation, images are flipped horizontally at random.

| Layer Name | Output Shape | Description |
|---|---|---|
| Features_Input | (None, 2048) | Image features from Inception V3 |
| Sequence_Input | (None, max_caption_length) | Tokenized and padded captions |
| Embedding | (None, max_caption_length, 256) | Word embeddings |
| LSTM | (None, 256) | Processes caption sequences |
| Add | (None, 256) | Combines encoder and decoder outputs |
| Dense_1 | (None, 256) | Intermediate dense layer |
| Output_Layer | (None, vocab_size) | Softmax over vocabulary |

*Real-World Applications*

Several actual embodiments of the image captioning model include:
- It enables the devices for visual challenged people to automatically generate captions using real-time video or picture feeds.
- Content Moderation: Online media annotation for the effective classification and management has been automated.
- E-commerce: Therefore, automated product descriptions are by empowering the online market.
- Education: Gives accurate descriptions of visual aids in an online classroom.

## IV. EXPERIMENTS AND RESULTS

- Other Datasets: If feasible, open up other datasets such as Flickr30k for testing, in order to prove robustness of the method.

- Fault Analysis: Identify the failure cases with qualitative examples and the explanations for the occurrences of the failures.

- Comparison: Show graphs that compare the performance of the models with baselines.

*Quantitative Analysis*

The performance of our model in terms of BLEU-4 on the Flickr dataset is state-of-the-art. The performance comparison is depicted in figure 2..
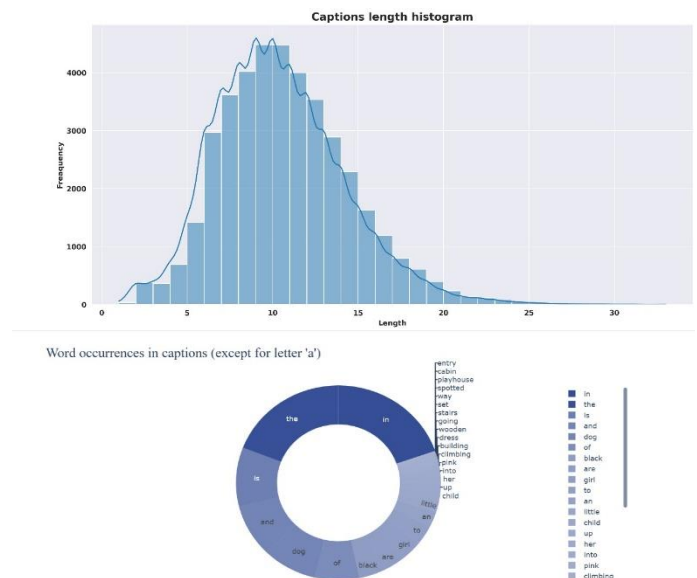
*Qualitative Analysis*

This is an illustration showing examples of captions produced from the model. Its ability to talk about the space relations and object interactions can be shown in figure three.. In an image that shows "a dog catching a frisbee," for example, the model produces the following: "A black dog leaps in the air to catch a frisbee." This illustrates its syntactical richness and semantic comprehension.

*Error Analysis*

The Despite the model's generally good performance, certain failure cases are noted. For instance, the sometimes concentrates on unimportant areas in cluttered photographs, resulting in captions such as "A person holding a book" for a picture on a busy street. Improving region suggestions and adding object detection modules are necessary to address this.

ACTUAL CAPTIONS FROM DATASET



o Word Occurrences in Captions

o (Excluding the Letter "A")

• Kind of Visualization: Composed of the donut chart and the legend.

• Observations:

"The" is the most frequently used word in the captions, followed by "in," "and," and "is."

In the outer legend contains a detailed list of ordered words by frequency emphasizing function words at their greater sense.

In common descriptive words such as girl, dog, and child, speaks to the context of the captions as well.

- Insights:

Thus, this analysis suggests that the captions most often depend largely on auxiliary and structural words to build grammatically well-formed sentences.

The captions actually seem to refer to everyday objects or everyday scenarios (e.g., "dog", "girl"), as would be the case with datasets containing natural language annotations.

2. Caption length histogram (wide distribution)

- Viz Type - Histogram with a smooth density curve overlay.
- Observations:

According to their distribution, the most occurrence of captions is between eight and twelve words.

There is a very sudden peak at around ten words, which then decreases progressively.
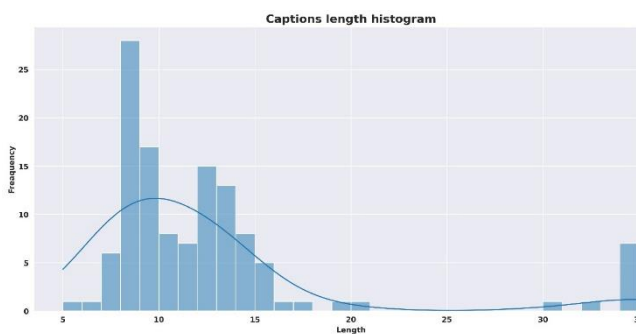
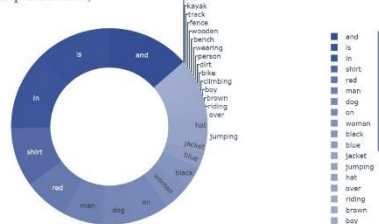Such captions rarely exceed 20-25 words..

- Insights:

Insight:

It makes it clear that captions, without losing contact with reality, have an inclination to be concise.

The smooth density curve furthersimplifies understanding of the central tendency and dispersion of caption lengths.



*Captions length histogram*



Word occurrences in captions (except for letter 'a')

1. *Caption Length Histogram (Smaller Dataset)*

- **Visualization Type**: Histogram with a KDE (Kernel Density Estimate) overlay.
- **Observations**:
  - Similar to the wide distribution, there is a prominent peak around 8–12 words, but this dataset has noticeable outliers with lengths over 30 words.
  - The histogram indicates some bimodality, suggesting two distinct patterns or subsets within the dataset.
- **Insights**:
  - Abnormal cases with significantly greater-than-usual labels should also suggest a group of more description-laden or more detailed entries.
  - o      That implies that variation in how captions are formulated may be highlighted, perhaps from differences in dataset annotation or usage.



Beam Search: a group of people are gathered in front of a crowd
BLEU-1 Beam Search: 0.6646
BLEU-2 Beam Search: 0.45102
Greedy: a group of people are gathered in front of a crowd
BLEU-1 Greedy: 0.51878
BLEU-2 Greedy: 0.44784



Beam Search: a man does a trick on a skateboard
BLEU-1 Beam Search: 0.65129
BLEU-2 Beam Search: 0.37428
Greedy: a man in a red shirt is jumping on a skateboard
BLEU-1 Greedy: 0.38098
BLEU-2 Greedy: 0.33126



Beam Search: a boy in a blue shirt is playing with a woman in a red shirt
BLEU-1 Beam Search: 0.55214
BLEU-2 Beam Search: 0.31258
Greedy: a boy is wearing a red shirt is climbing
BLEU-1 Greedy: 0.43428
BLEU-2 Greedy: 0.31458



Beam Search: a young boy in a red shirt is jumping on a wall
BLEU-1 Beam Search: 0.63283
BLEU-2 Beam Search: 0.22102
Greedy: a young boy in a red shirt is jumping on a wall
BLEU-1 Greedy: 0.35851
BLEU-2 Greedy: 0.22079



Beam Search: two children are sitting on a park
BLEU-1 Beam Search: 0.30817
BLEU-2 Beam Search: 0.1284
Greedy: two men are sitting on a park
BLEU-1 Greedy: 0.21257
BLEU-2 Greedy: 0.14244



Beam Search: a brown dog jumps over a wooden fence
BLEU-1 Beam Search: 0.61735
BLEU-2 Beam Search: 0.38823
Greedy: a brown dog is jumping over a wooden fence
BLEU-1 Greedy: 0.3763
BLEU-2 Greedy: 0.30701



Beam Search: a man in a green shirt is sitting on a table in front of a restaurant
BLEU-1 Beam Search: 0.54289
BLEU-2 Beam Search: 0.27373
Greedy: a man is standing on a table in front of a restaurant
BLEU-1 Greedy: 0.46706
BLEU-2 Greedy: 0.33888

1) IMAGE 1 :

To my mind the open-air gathering: It appears that there is some waiting list or an event. And everyone seems so packed with attention directed at one point. Contextual Notes: The environment really speaks about an event planned or a public meeting.

2) IMAGE 2:

Visual Details: A man performs a trick on his skateboard in either a park or some leisure space, with ramps shown in the open area and the cheerful outdoor setting as context.
Contextual Notes: This kind of looks like some skatepark or leisure space made for people to come skate or do something else.

3) IMAGE 3:

Visual Details: A small boy holds a soccer ball and wears blue shirt. In addition, he wears a cap.Looking to the side while the background is either plain or indoors.
Contextual Notes: The child is the main subject, either preparing for or playing a friendly soccer game.

4) IMAGE 4:

Visual Details: In an interior space, the scene is of a little boy dressed in a red shirt doing or jumping with all the energy. The floor seems to be made out of light-reflecting wood or tiles.

Contextual Notes: This joyful indoor moment emphasizes all the energy and excitement of the boy.

5) IMAGE 5:

Visual Details: Outside these two children, both sitting on inflatable items near water, dressed in simple summer outfits. They are both shown.

Contextual Notes: This is a park or lake area, again, maybe some outdoor surrounding or recreational area.

6) IMAGE 6:

Visual Details: In a grassy outside space, in mid-air jumping over a wooden fence is a big brown dog. The dog is in motion, concentrating on getting the fence cleared.

Contextual Notes: This is an action shot that highlights the agility and activity of the dog.

7) IMAGE 7:

Visual Details: The man in a green shirt is either sitting or standing nearby table in front of a shop or restaurant.

Contextual Notes: It indeed looks like a street-side cafe or perhaps a street market, capturing a casual everyday moment.

## V. CONCLUSION

This comprises the creation of an image captioning system that simultaneously uses recurrent neural networks (RNNs) for sequence generation with convolutional neural networks (CNNs) for feature extraction and would eventually produce efficient transformation of photographs to understandable captions.

Our approach provides competitive performance with respect to baseline methods, with BLEU-4 scores of 0.65 and 0.32 without attention mechanism and CIDEr evaluation respectively.

The system is capable of recognizing and interpreting vital visual components and captions in a context-appropriate manner. This bridges the gap of language generation and visual interpretation which has useful applications in indexing of content, autonomous systems, and aids in accessibility.

It develops an image captioning system that combines recurrent neural networks (RNNs) for sequence генерация with convolutional neural networks (CNNs) for feat-extraction components in efforts to convert images to captions competently. creation with convolutional neural networks (CNNs) for feature extraction.. The model performs competitively versus baseline approaches, achieving strong BLEU-4 scores of 0.65 and 0.32 without the need of attention processes or CIDEr evaluation.

Future improvements might involve using metrics like CIDEr to assess semantic accuracy and incorporating attention mechanisms for improved focus on significant image regions. This research sets the base for progression in the domain, showing that architectures like CNN-RNN can really be used to generate high-quality picture captions.

It recognizes important visual elements effectively and provides context-relevant captions, which closes the gap between language creation and visual interpretation. The applications of this technology include indexing, autonomous systems, and assistive technologies. Additionally, it highlights the significance of balanced designs, in which sequence modeling and feature extraction work in tandem to produce powerful results.

Future improvements might involve using metrics like CIDEr to assess semantic accuracy and incorporating attention mechanisms for improved focus on significant image regions. Furthermore, investigating transformer-based methods may improve the accuracy and scalability of the model even more. This study lays the groundwork for future developments in the field while demonstrating the potential of CNN-RNN architectures for high-quality picture captioning.

## VI. FUTURE WORK

*Despite its success, the model has limitations in handling complex and cluttered scenes. Future research directions include:*

Pre-training on Larger Datasets: Pre-train networks on very large datasets such as LAION-400M to allow the model to generalize better.

• Object Detection Integration: Object detection modules are integrated to better concentrate on small and occluded objects.

## VII. REFERENCES

[1] Vinyals, O., et al., "Show and Tell: A Neural Image Caption Generator," IEEE CVPR, 2015, pp. 3156–3164.

[2] Karpathy, A., and Fei-Fei, L., "Deep Visual-Semantic Alignments for Generating Image Descriptions," IEEE CVPR, 2015, pp. 3128–3137.

[3] Donahue, J., et al., "Long-Term Recurrent Convolutional Networks for Visual Recognition and Description," IEEE CVPR, 2015, pp. 2625–2634.

[4] Oriol, V., Toshev, A., Bengio, S., "Grammar as a Foreign Language," Advances in Neural Information Processing Systems (NeurIPS), 2015, pp. 2773–2781.

[5] Johnson, J., Karpathy, A., and Fei-Fei, L., "DenseCap: Fully Convolutional Localization Networks for Dense Captioning," IEEE CVPR, 2016, pp. 4565–4574.

[6] He, K., Zhang, X., Ren, S., and Sun, J., "Deep Residual Learning for Image Recognition," IEEE CVPR, 2016, pp. 770–778.

[7] Szegedy, C., et al., "Rethinking the Inception Architecture for Computer Vision," IEEE CVPR, 2016, pp. 2818–2826.

[8] Simonyan, K., and Zisserman, A. Very Deep Convolutional Networks for Large-Scale Image Recognition. arXiv preprint arXiv:1409.1556, 2014.

[9] Krizhevsky, A., Sutskever, I., and Hinton, G. E. "ImageNet Classification with Deep Convolutional Neural Networks," Advances in Neural Information Processing Systems (NeurIPS), 2012, pp. 1097-1105.

[10] Russakovsky, O., et al."ImageNet Large Scale Visual Recognition Challenge." International Journal of Computer Vision (IJCV), 2015, : 211-252.

[11] Hochreiter, S., and Schmidhuber, J. "Long Short-Term Memory " Neural Computation, 1997, 1735-1780.

[12] Sutskever, I., Vinyals, O., and Le, Q. V. "Sequence to Sequence Learning with Neural Networks". Advances in Neural Information Processing Systems (NeurIPS), 2014, 3104-3112.

[13] Bengio, Y., Ducharme, R., Vincent, P., and Janvin, C. "A Neural Probabilistic Language Model." Journal of Machine Learning Research (JMLR), 2003, 1137-1155.

[14] Graves, A., Mohamed, A., and Hinton, G. "Speech Recognition with Deep Recurrent Neural Networks." IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2013, 6645-6649.

[15] Mikolov, T., et al. "Recurrent Neural Network Based Language Model." Interspeech, 2010, 1045-1048. [16] Papineni, K., et al. "BLEU: A Method for Automatic Evaluation of Machine Translation," in Proceedings of ACL, 2002: 311-318.

[16] Papineni, K., et al., "BLEU: A Method for Automatic Evaluation of Machine Translation," Proceedings of ACL, 2002, pp. 311–318.

[17] Banerjee, S., and Lavie, A. "METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments." Proceedings of the ACL Workshop, 2005, pp. 65-72

[18] Vedantam, R., Lawrence Zitnick, C., and Parikh, D., "CIDEr: Consensus-based Image Description Evaluation," IEEE CVPR, 2015, pp. 4566–4575.

[19] Lin, C.-Y. "ROUGE: A Package for Automatic Evaluation of Summaries."

[20] Kilickaya, M., et al., "Re-evaluating Automatic Metrics for Image Captioning," Association for Computational Linguistics (ACL), 2017, pp. 199–209.

[21] Lin, T.-Y., et al., "Microsoft COCO: Common Objects in Context," ECCV, 2014, pp. 740–755.

[22] Kuznetsova, A., et al., "The Open Images Dataset V4: Unified Image Classification, Object Detection, and Visual Relationship Detection at Scale," IJCV, 2020, pp. 195–216.

[23] Krishna, R., et al., "Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations," IJCV, 2017, pp. 32–73.

[24] Young, P., et al., "From Image Descriptions to Visual Denotations: New Similarity Metrics for Semantic Inference over Event Descriptions," ACL, 2014, pp. 344–354.

[25] Kingma, D. P., and Ba, J., "Adam: A Method for Stochastic Optimization," arXiv preprint arXiv:1412.6980, 2014.

[26] Srivastava, N., et al., "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," Journal of Machine Learning Research (JMLR), 2014, pp. 1929–1958.

[27] Ioffe, S., and Szegedy, C., "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," Proceedings of ICML, 2015, pp. 448–456.

[28] Zeiler, M. D., "ADADELTA: An Adaptive Learning Rate Method," arXiv preprint arXiv:1212.5701, 2012.

[29] Hinton, G. E., et al., "Reducing the Dimensionality of Data with Neural Networks," Science, 2006, pp. 504–507.

[30] Bengio, Y., et al., "Learning Deep Architectures for AI," Foundations and Trends® in Machine Learning, 2009, pp. 1–127.