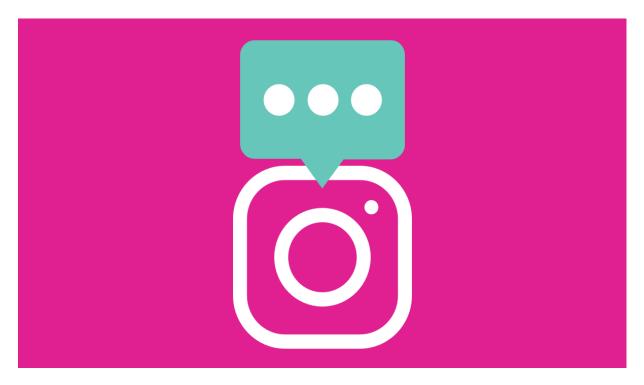
Project Proposal



Automatic Image Captioning

Group 16	
Project 1 Title	Automatic Image Captioning
Team Members	Achanta Ajitha
	Chandan Kumar
	Jammi Venkata Sai Harsha
Mentor	Brahmani
Date	11-June-2023

Objective:

We will develop models which will automatically detect the accurate captions for images. We will use different deep neural networks to detection an image. The results will be combined and used by a captioning module that generates image captions through natural language processing techniques. We will deploy the model by building a dashboard.

Business Use Cases:

Automatic image captioning has several potential business use cases across different industries. Here are a few examples:

- **1. E-commerce:** Automatic image captioning can be used to generate accurate and descriptive captions for product images. This can improve the searchability of products, enhance the user experience, and provide more information to potential customers. For example, a clothing retailer can automatically generate captions that include details like the type of garment, color, size, and style.
- 2. Social Media and Marketing: Image captioning can help businesses automatically generate captions for their social media posts, saving time and effort for content creators. This can be particularly useful for platforms like Instagram, where visuals play a significant role. Automatic image captioning can also aid in creating engaging and relevant hashtags based on the content of the image.
- **3. Content Moderation:** Image captioning can assist in content moderation by automatically analyzing images and generating captions that describe their content. This can be used to identify and flag inappropriate or sensitive content, ensuring a safer online environment.
- **4. Media and Publishing:** In the media and publishing industry, automatic image captioning can be used to add captions to news articles, blogs, or photo galleries. This can provide additional context to the images and improve the overall storytelling experience.
- **5. Accessibility:** Automatic image captioning can benefit individuals with visual impairments by providing them with textual descriptions of images. This can make online content more inclusive and enable visually impaired users to access and understand visual information.
- **6. Surveillance and Security:** Image captioning can be used in security systems to automatically analyze and describe images captured by surveillance cameras. This can aid in identifying and flagging suspicious activities or objects, enhancing overall security measures.

Problem Statement:

The problem of automatic image captioning aims to develop a system that can accurately generate descriptive and contextually relevant captions for images. The challenge lies in bridging the gap between visual information captured in an image and its textual representation, requiring the model to understand and interpret the visual content effectively.

Key Challenges:

1. Image Understanding: Developing a model that can accurately understand the visual content of an image, including objects, scenes, relationships, and context, is a fundamental challenge. The model should possess the ability to extract meaningful features and comprehend the semantic information present in the image.

- **2. Natural Language Generation:** Generating captions that are grammatically correct, coherent, and semantically meaningful is a significant challenge. The model needs to understand the appropriate language style, structure, and vocabulary to produce captions that effectively describe the image content.
- **3. Scalability and Efficiency:** Developing an automatic image captioning system that is scalable and efficient is essential for real-world applications. The model should be able to process images and generate captions in a timely manner, even when dealing with large-scale datasets or real-time scenarios.

Literature Review:

In this section, we discuss the three main categories of existing image captioning methods: template-based image captioning, retrieval-based image captioning, and novel caption generation. Template-based techniques have fixed templates with blank slots to generate captions. In these systems, the different objects, actions and attributes are first identified and then the gaps in the templates are filled. For example, Farhadi et al. [1] use three different elements of a scene to fill the template slots for generating image captions. A Conditional Random Field (CRF) is leveraged by Kulkarni et al. [2] to detect the objects, attributes, and prepositions before filling in the blanks. Template-based approaches are able to generate grammatically correct captions, but since the templates are predefined, it cannot generate variable-length captions.

In retrieval-based methods, captions are retrieved from a pool of existing captions. Retrieval based methods initially find images that are visually similar images to the given image along with their captions from the training data set. These captions are called candidate captions. The captions for the given image are selected from this caption set [3], [4]. These types of systems produce general and grammatically correct captions. However, they cannot generate more descriptive and semantically correct captions.

Novel captions can be generated from visual and multimodal spaces. In these types of systems, the visual content of the image is first analysed and then captions are generated from the visual content using a language model [5], [6], [7], [8]. These approaches can generate new, more semantically accurate captions for each image. Most novel caption generation techniques employ deep machine learning. Therefore, in this paper we focus primarily on deep learning based novel image caption generating methods.

Data and Data Pre-Processing:

Data and data pre-processing are crucial components in developing an automatic image captioning system. Data pre-processing is a critical step that ensures the dataset is in a suitable format for training the image captioning model. It helps to normalize, clean, and transform the data, making it more manageable and compatible with the model architecture and training process.

Here's an overview of the data requirements and pre-processing steps involved in our project:

1. Image Data: We need a large dataset of paired images and their corresponding captions for training. These images can be sourced from various publicly available image datasets such as Flickr30k or Flickr8k. The dataset should cover a wide range of subjects, scenes, and objects to ensure diversity and generalization of the model.

- **2. Caption Data:** Each image in the dataset will have one or more human-generated captions associated with it. The captions will accurately describe the content of the image and provide sufficient context. Ensuring high-quality, relevant, and descriptive captions is essential for training an effective image captioning model.
- **3. Text Pre-processing:** Before training the model, the caption texts undergo pre-processing steps, including:
 - Tokenization: The captions are split into individual words or tokens, creating a vocabulary.
 - Lowercasing: Converting all words to lowercase to ensure consistent word representations.
 - Removing Punctuation: Punctuation marks and special characters are often removed to simplify the text.
 - Removing Stop Words: Commonly used stop words (example, "the," "a," "and") that carry little semantic value are removed.
- **4. Image Pre-processing:** Images need to undergo pre-processing steps to ensure consistency and compatibility with the model. Common pre-processing techniques include:
 - Resizing: Resizing the images to a fixed dimension to ensure uniformity. Common size includes 224x224x3 pixels.
 - Normalization: Normalizing the pixel values to a common scale, typically between 0 and 1, to facilitate efficient model training.
 - Data Augmentation: Applying random transformations such as rotations, flips, or crops to augment the dataset and improve model generalization.
- **5. Data Split:** The dataset is typically divided into three subsets: training, validation, and testing. The training set is used to train the model, the validation set helps in tuning hyperparameters and monitoring performance, while the testing set is used to evaluate the final model's performance.
- **6. Word Indexing:** Each unique word in the pre-processed captions is assigned a unique index. This indexing enables mapping between words and their numerical representations, facilitating model training and inference.

Model and Implementation:

Automatic image captioning involves generating textual descriptions for images using deep learning techniques. One popular approach for implementing automatic image captioning is to use a combination of Convolutional Neural Networks (CNNs) for image feature extraction and Recurrent Neural Networks (RNNs) for generating captions.

Here is a high-level overview of a typical model and implementation for automatic image captioning:

1. Dataset Preparation:

 Gather a large <u>Flickr30k</u> or <u>Flickr8k</u> dataset of images paired with corresponding captions. These captions serve as the ground truth for training the model.

- Preprocess the images by resizing them to a fixed size and apply necessary normalization or augmentation techniques.
- Tokenize the captions into individual words or subwords to create a vocabulary.

2. CNN Feature Extraction:

- We will use a pre-trained CNN model such as VGG16 to extract high-level features from the input images.
- Remove the last fully connected layers of the CNN to obtain a fixed-length feature vector representing the image content.
- These features will capture the visual information necessary for generating captions.

3. RNN Caption Generation:

- Initialize RNN such as a Long Short-Term Memory (LSTM) as the caption generator.
- The image features from the CNN will use as the initial hidden state of the RNN.
- Feed the caption sequence into the RNN one word at a time, predicting the next word in the sequence based on the previous words and the image features.
- Train the RNN using a loss function such as cross-entropy loss, comparing the predicted captions with the ground truth captions.

4. Model Training:

- Combine the CNN feature extraction and RNN caption generation components into an end-to-end model.
- Initialize the parameters of the model with random values or will use pre-trained weights for the CNN part.
- Use the prepared dataset to train the model, optimizing the parameters using backpropagation and gradient descent.
- Adjust hyperparameters such as learning rate, batch size, and regularization techniques, to improve performance and prevent overfitting.

5. Caption Generation:

- During inference provide an input image to the trained model.
- Pass the image through CNN to extract the image features.
- Initialize the RNN with the extracted features as the initial hidden state.
- Generate captions by feeding the output of the RNN back into itself, iteratively
 predicting the next word in the sequence until an end token is generated, or a
 maximum caption length is reached.
- Apply techniques like sampling to improve the diversity and quality of generated captions.

6. Evaluation and Fine-tuning:

- Evaluate the quality of the generated captions using any metrics like BLEU, METEOR, or CIDEr.
- Fine-tune the model based on the evaluation results to improve performance.
- Repeat the training and evaluation process until satisfactory results are achieved.

Conclusion:

The above steps provide a high-level overview, and there can be variations and improvements to this approach. There may be modifications to the architecture, incorporate attention mechanisms, or experiment with different training strategies to enhance the captioning performance.

References:

[1 William Fedus, Ian Goodfellow, and Andrew M Dai. Maskgan: Better text generation. arXiv preprint arXiv:1801.07736, 47, 2018. https://arxiv.org/abs/1801.07736

- [2] Girish Kulkarni, Visruth Premraj, Sagnik Dhar, Siming Li, Yejin Choi, Alexander C Berg, and Tamara L Berg. Baby talk: Understanding and generating image descriptions. IEEE Transactions on Pattern Analysis and Machine Intelligence, 35:2891–2903, June 2013. http://acberg.com/papers/baby_talk.pdf
- [3] Yunchao Gong, Liwei Wang, Micah Hodosh, Julia Hockenmaier, and Svetlana Lazebnik. Improving image-sentence embeddings using large weakly annotated photo collections. European Conference on Computer Vision. Springer, pages 529–545, 2014. https://link.springer.com/chapter/10.1007/978-3-319-10593-2 35
- [4] Peter Young Micah Hodosh and Julia Hockenmaier. Framing image description as a ranking task: Data, models and evaluation metrics. Journal of Artificial Intelligence Research, 47:853–899, 2013.

https://www.ijcai.org/Proceedings/15/Papers/593.pdf

- [5] Ryan Kiros, Ruslan Salakhutdinov, and Richard S Zemel. Unifying visual-semantic embeddings with multimodal neural language models. Workshop on Neural Information Processing Systems (NIPS), 2014. https://arxiv.org/abs/1411.2539
- [6] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. International Conference on Machine Learning, 2048- 2057, 2015. https://proceedings.mlr.press/v37/xuc15.html
- [7] Ting Yao, Yingwei Pan, Yehao Li, Zhaofan Qiu, and Tao Mei. Boosting image captioning with attributes. IEEE International Conference on Computer Vision (ICCV), pages 4904–4912, 2017.

https://arxiv.org/abs/1611.01646

- [8] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. Image captioning with semantic attention. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 4651–4659, 2016. https://arxiv.org/abs/1603.03925
- [9] https://towardsdatascience.com/a-guide-to-image-captioning-e9fd5517f350
- [10] https://github.com/Jasminehh/automatic image captioning
- [11] https://fairyonice.github.io/Develop_an_image_captioning_deep_learning_model_using_Flickr 8K data.html