

PYTHON-BASED IMAGE CAPTIONING USING CNN AND LSTM

Abstract

Our most important sense is vision. The capacity to see has been used by software engineers to create more dynamic, intelligent, and easily accessible software through visuals. There are circumstances, nevertheless, in which an image might not be enough. Alternative text may be provided to avoid bandwidth limits and offer a more accessible experience if further context is required. The manual explanation falls short in an age where there is simply a great deal of photographs to describe. Deep learning can integrate image processing and natural language processing, allowing computers to independently provide explanations for images. This service can be provided through a user-friendly web interface, where customers can simply upload the photos, they want to be described. This makes it possible for anybody to easily take advantage of the capabilities of this deep learning technique and adaptive image descriptor capability through a simple API, with the computationally demanding chores being abstracted away.

Keyword: The manual explanation falls short in an age where there is simply a great deal of photographs to describe.

Authors

Dr. A. Suneetha

Associate professor
Department of CSE
KKR & KSR Institute of Technology &
Sciences
Guntur, Andhra Pradesh.

Dr. K. Ratna Babu

Lecturer in Computer Engineering
Department
Government Polytechnic
Addanki.

I. INTRODUCTION

The advent of electronic computers and those with stored programs laid the foundation for the exploration of intelligent systems. This sparked the question fueled by human curiosity: "Can a machine emulate human thinking and behavior, leveraging the computational power of computer systems?" This inquiry led to the development of Intelligent Systems, with the aspiration of imbuing robots with a level of intelligence comparable to that highly esteemed in humans. After decades of development, artificial intelligence (AI) has become a tangible reality. It has begun assuming tasks beyond human capacity. In contrast to the innate intelligence displayed by humans and other animals, artificial intelligence, as defined by Wikipedia, pertains to the intelligence demonstrated by computers. In computer science, the examination of "intelligent agents" encompasses any machine capable of perceiving its environment and taking actions to enhance the likelihood of achieving its objectives.

1. Related Work: They propose adopting a modified approach involving the reinstatement of the encoder Network, previously employed in a deep CNN. CNNs have consistently demonstrated their ability to present images effectively for various vision applications by converting input pictures into fixed-length arrays. Consequently, the standard practice involves initially training a convolutional neural network for image classification tasks and subsequently using the final hidden layer as input for the RNN decoder, which generates words, effectively employing the network as an "encoder" for images. This model, named Neural Image Caption (NIC), has proven to be a reliable solution to the problem. Notably, the neural network can be fully trained using stochastic gradient descent. Additionally, the model incorporates the latest language and vision sub-networks, often pre-trained on larger datasets, thus benefiting from more extensive data. Lastly, it outperforms contemporary methodologies, yielding superior results.

Images can be analyzed using a blend of probabilistic and neural methodologies. Recent progress in statistical machine translation (MT) suggests that with a well-constructed sequence model, we can attain positive results by actively boosting the likelihood of precise translation through an "end-to-end" method, applied for both training and prediction. These models leverage a Recurrent Neural Network (RNN) to convert the input, which varies in length, into a fixed-dimensional vector, ultimately producing a predefined output. Therefore, it is customary to apply a comparable "conversion" procedure when elucidating an image, rather than a sentence in the source language.

To tackle the challenge of sequence generation and translation, which stands as the primary hurdle in RNN development and training, a specialized form of recurrent network known as LSTM was devised and implemented with remarkable success. The pivotal feature of LSTM memory cells involves encoding information at each time step based on inputs observed up to the present moment. These "gates" serve to regulate the behavior of the recurrently incorporated cells or layers. Depending on whether the gate is set to one (1) or zero (0), the gated layer may preserve its value. Specifically, there are

three gates—the input gate, the forget gate, and the output gate—that determine whether the system should read its input, disregard the current cell value, or produce an output.

- 2. Existing System:** Humans are inherently the foremost and most adept individuals for interpreting an image, yet with the staggering volume of photos captured and shared daily, there arises a question: who will undertake this task? This is where human potential and vigor reach their zenith. When a person approaches their limits, technology steps in to provide a solution. One such technique is Deep Learning (DL), a pivotal facet of AI. In contemporary systems like "show and tell," the ANN, emulating human brain behavior, has been applied with notable success. Convolutional neural networks (CNNs) are employed to emulate human-like functions, particularly in tasks like image description. While Convolutional neural networks are frequently harnessed for image processing, owing to their capacity for high precision, recurrent neural networks (RNNs) excel in handling textual and auditory information, given their adeptness at swiftly processing sequential data.

In both the Bottom-Up and Top-Down approaches, it is imperative to pose a question to the system for it to generate a suitable description. This interactive process of inquiry and response formation is known as Visual Question Answering (VQA). The system needs to be autonomous and self-reliant, as it would be impractical for individuals to continuously pose questions about every single image to obtain comprehensive descriptions.

- 3. Proposed System:** Upon careful examination of the issue and the current system's approach, we discerned an opportunity for innovative problem-solving that could yield superior outcomes. The preceding methods employed to address this challenge were pioneering. Notably, "Show and Tell" astutely harnessed the advancements in transfer learning within image processing, distinguishing itself from contemporaries who viewed the issue solely through a linguistic lens. This approach was indeed apt: Image captioning encompasses elements of both language processing and image processing. By systematically focusing on each of these domains independently, they were able to leverage the most effective strategies within each. The integration of these dual tactics proved highly successful. In their approach, CNN was employed for image processing to generate input for the subsequent Recurrent Neural Network (RNN)-based text generation. Extensive research and analysis have solidified CNNs as leaders in addressing visual challenges. Given that visual content inherently involves images, the CNN integrated into their methodology demonstrated exceptional performance, further enhanced by human adjustments that enabled the RNN to provide more refined captions.

In pursuit of a self-triggering system capable of producing high-quality captions, we chose to integrate an advanced deep learning technique: Generative Adversarial Networks (GANs). This architecture revolves around a dual learning process that encompasses both a generator and a discriminator. The generator's function is to handle the given input and generate the desired system output. In parallel, the discriminator assumes the role of an anticipatory assessor of the output, mimicking a user before granting direct access to the generated result. Leveraging the available information, the discriminator evaluates the adequacy of the output generated by the system. When the discriminator deems the output satisfactory, it signifies.

As the model undergoes training over several epochs, the generator within the augmented GAN model refines its ability to produce output, driven by the discriminator's feedback for continuous improvement. Simultaneously, the discriminator becomes increasingly proficient in discerning subpar output. This dual refinement process is a pivotal factor influencing GANs' capacity to yield high-quality outcomes in a machine-learning context. By amalgamating the strengths of previous approaches into a unified system, a novel proposition emerges. This new system builds upon its predecessor and progresses to generate improved results by addressing the limitations of existing systems through its evolutionary journey.

II. METHODOLOGY

1. Problem Scope: The problem-solving methodologies employed in the past were characterized by their innovative approaches. While some approached the challenge from a linguistic standpoint, the Show and Tell method introduced a brilliant solution by capitalizing on the advancements in transfer learning within image processing. This insight was astute: captioning images encompasses both natural language processing and image processing domains. By meticulously addressing each facet individually and subsequently integrating them, they were able to leverage the most potent techniques from both disciplines, resulting in remarkable success. Their pioneering approach involved the implementation of a state-of-the-art CNN for image processing, serving as the input framework for an RNN utilized in text synthesis. Through rigorous research and thorough analysis, they honed CNNs to excel in tackling visual challenges, establishing them as pioneers in their field. This specialized CNN exhibited exceptional performance, capitalizing on the visual richness inherent to media, thereby enhancing the RNN's proficiency in generating superior captions.

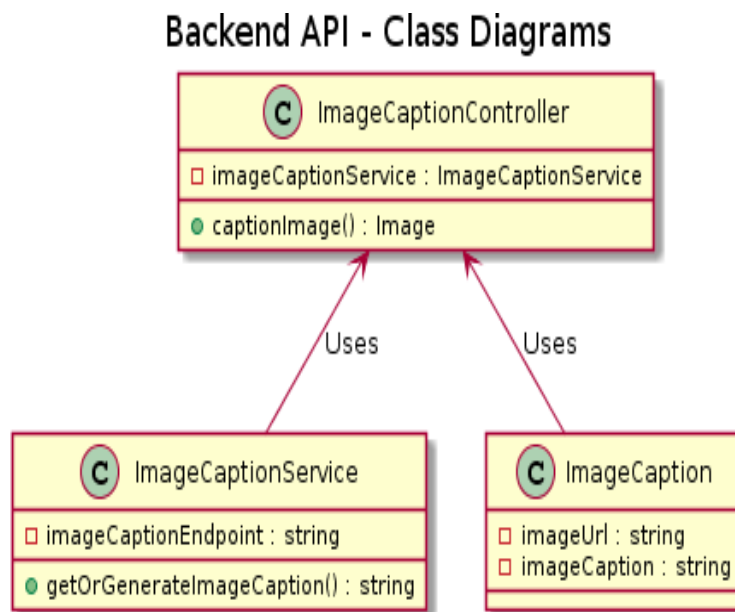


Figure 1: Backend API

Designed with dependency injection (DI) at its core, the backend API prioritizes

loose coupling among its components. This implies that its dependents can operate without worrying about its internal dependencies, as long as the components conform to the established interface or contract. This DI-based approach also facilitates unit testing by allowing for the simulation of component dependencies, enabling isolation and observation during testing.

2. ML Module :

Machine Learning Module - Class Diagrams

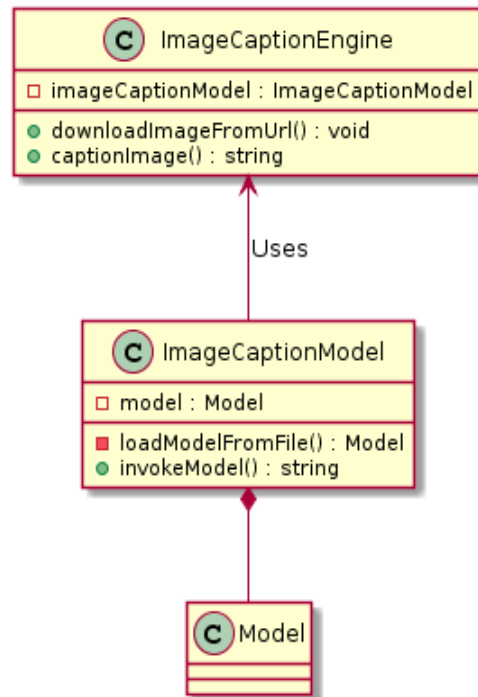


Figure 2: ML module

The module tasked with generating captions for photographs encompasses the machine learning component. Central to this process is the machine learning model we designed specifically for captioning images. This model undergoes rigorous training across multiple epochs or iterations, a pivotal phase in the process of operationalizing the machine learning model for practical use.

Once the training process concludes, the model is stored on disk, obviating the necessity of reiterating the entire procedure whenever a caption is needed for an image. Upon receiving a request for an image caption, the model is loaded from the disk and employed for inference. This methodology substantially enhances the processing speed of the model. In this context, "model" pertains to the pre-trained and preserved model.

3. Frontend:

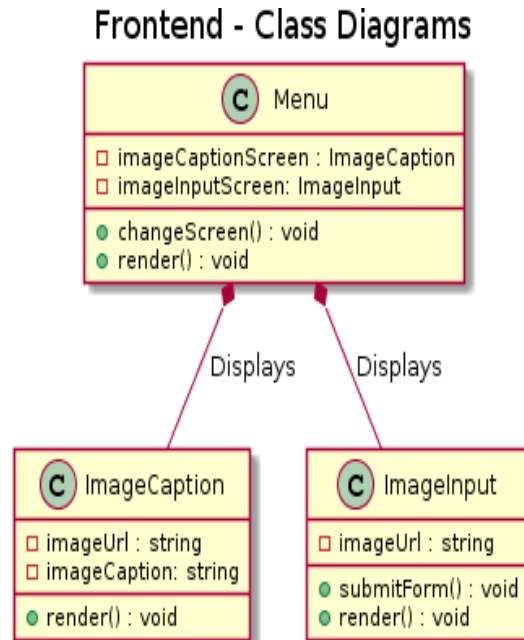


Figure 3: Front End

The majority of the classes in the frontend are React functional components, which take care of handling user interface (UI) display and HTTP request communication with the backend API. The user's contact with the system began with the Menu class. It manages the user interaction flow and is in charge of altering the visible screen to reflect the user's selected course of action. To do this, it keeps track of the Image Caption and Image Input screens, calls respective render methods when the screens are modified, and injects data into both of them. All of these classes have the render function since it is how React recognizes them for screen display. The shift Displaying the required information is handled by the scree method.

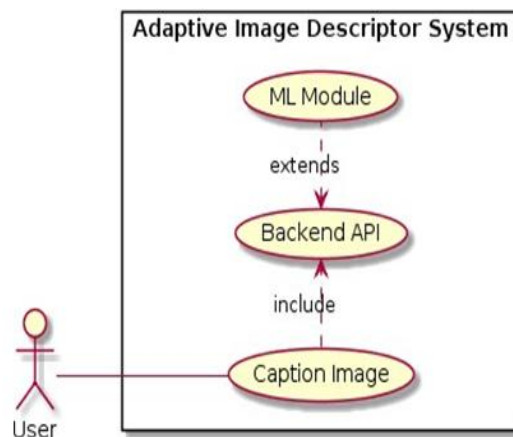


Figure 4: Captioning Image process

III. ARCHITECTURE

The system's components' non-linear interactions with one another are a defining feature of asynchronous architectures. This architecture's elements are also referred to as services. Services use asynchronous methods of communication with one another; the majority of the time, they are informed that their request has been approved rather than receiving a prompt answer. The response is stored in a shared database or message queue when the request has been fully processed, and it is then retrieved. One of the most popular ways of communicating in this architecture is through message queues. The messages are transferred and stored in a message queue. On the message queue, services may subscribe to a channel and either push messages there or retrieve messages that have already been pushed there. A channel serves as a point of contact for the communication between several services.

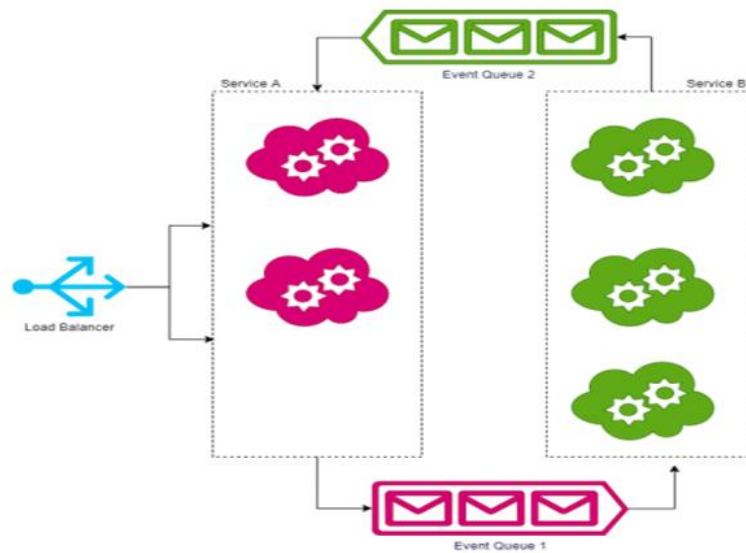


Figure 5: Process Architecture

1. Specific Architecture

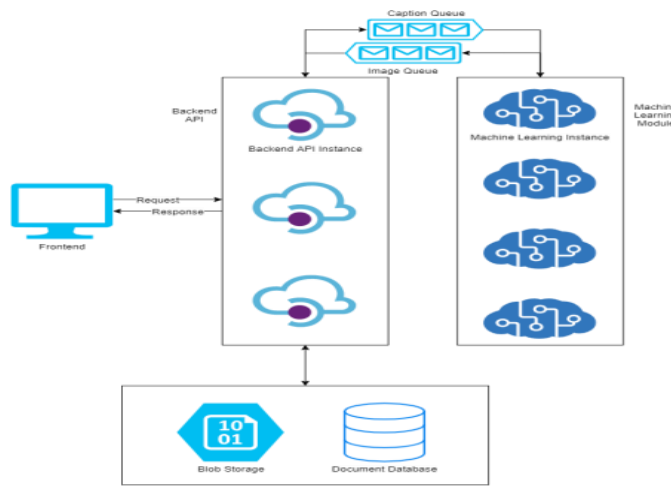


Figure 6: Specific Architecture



Figure 7: Object recognition (left) and object detection (right).

Utilizing Deep Learning (DL) techniques for Object Identification has become a widely adopted approach in object recognition. Convolutional neural networks (CNN) and other DL models facilitate the identification of objects by autonomously discerning their distinctive features. For instance, CNN can become proficient at distinguishing between cats and dogs through exposure to an extensive set of training images, thus discerning the unique attributes that differentiate them. This process necessitates a substantial labeled dataset and a network architecture capable of learning these features and constructing the model from scratch. While it demands a significant amount of training data and the alignment of the CNN's layers and weights, the outcomes can be remarkably impressive.

IV. RESULTS AND DISCUSSION

The file that defines all the libraries and files needed and utilized by the project is

called `a.csproj`. The backend makes use of the NuGet packages and the features they offer. Advanced JSON capability is offered by the `Newtonsoft.Json` packages. The project uses it to deliver data serialization and deserialization in JSON format. Deserializing JSON data into POCO (Plain Old C# Object) classes and serializing POCO classes back into JSON data are two different processes. This is due to the backend's JSON-based API being available. To avoid several allocations, utilize the Object Pool extension. Although the CLR, which supports ASP.NET Core, is effective at executing garbage collection, preventing the creation of superfluous objects is one of the greatest methods to enhance the speed of the application. We may define a pool of once-created objects using the Object Pool package. Every time a certain item is needed, it is removed from the pool. The item is returned to the pool once it has served its purpose. We avoid making additional allocations this way. Our preferred database, MongoDB, may be reached using the MongoDB Driver package. We may create a connection with the database by using the driver. The driver assists us in getting a hold of the database collections after the connection has been made. Operations that modify the collection are permitted with this handle.

The major backend entry point is the Startup class. All of the services that are used are listed and set up here. All of the backend's services must be configured using the `Configure Services` method. The `Configure` method manages the activities that must be taken by the system setup. Additionally, the singletons used by other classes in the system are initialized by this class. It uses the Dependency Injection package to perform dependency injection. The Startup class must initialize before the backend can begin processing requests.

V. CONCLUSION

While developing this method, we built a model for picture captioning from scratch and demonstrated how deep learning can be used to create accurate captions in languages like English. We used the Flickr8k dataset to develop our algorithm, and finally, we were able to create captions with a mediocre level of accuracy. We also concluded that the accuracy would increase with a larger dataset while the losses would decrease. Automatic captioning of photos is a relatively new activity, but thanks to the research of academics in this area, great advancement has been made. We believe there is still a great deal of space to improve the efficacy of photo captioning. First off, the rapid advancement of deep neural networks will surely improve the efficacy of picture description creation by utilizing more effective network designs as language models and/or visual models. Second, since captions are collections of words whereas pictures are made up of things arranged in space, it is crucial to look at the presence and arrangement of visual concepts in captions. Investigating how to address issues with picture captioning in various unique instances may also be fascinating. Our study was conducted utilizing the Flickr8k dataset. In contrast, for future improvement, we may employ progressively bigger datasets, such as Flickr30k or MSCOCO, to produce more accurate captions.

REFERENCES

- [1] Wu Y, et al. (2016) Google's neural machine translation system: Bridging the gap between human and machine translation CoRR abs/1609.08144.
- [2] Meltzoff A, Kuhl P, Movellan J, Sejnowski T (2009) Foundations for a new science of learning. 325:284-8.
- [3] Hinton G, et al (2012) Deep neural networks for acoustic modeling in speech recognition: The shared view of four research groups. *IEEE Signal Processing Magazine* 29(6):82-97.
- [4] Mishra J, Saha I (2010) Artificial neural networks in hardware: A survey of two decades of progress.

- Neurocomputing* 74(1):239-255. Artificial Brains.
- [5] Maier HR, Dandy GC (2000) Neural networks for the prediction and forecastin of water resources variables: a review of modeling issues and applications. *Environmental Modeling and Software* 15(1):101-124.
- [6] Bhute AN, Meshram BB (2014) Text-based approach for indexing and retrieval of image and video: A review. *CoRR* abs/1404.1514.
- [7] Karpathy A, Fei-Fei L (2017) Deep visual-semantic alignments for generating image descriptions. *IEEE Trans. Pattern Anal. Mach. Intell.* 39(4):664-676.
- [8] Vinyals O, Toshev A, Bengio S, Erhan D (2014) Show and tell: A neural image caption generator. *CoRR* abs/1411.4555.
- [9] Lipton ZC, Kale DC, Elkan C, Wetzel RC (2015) Learning to diagnose with LSTM recurrent neural networks. *CoRR* abs/1511.03677.
- [10] Sharma, Grishma, Priyanka Kalena, Nishi Malde, Aromal Nair, and Saurabh Parkar. "Visual Image Caption Generator Using Deep Learning." Available at SSRN 3368837 (2019).
- [11] Tan, Ying Hua, and Chee Seng Chan. "Phrase-based image caption generator with hierarchical LSTM network." *Neurocomputing* 333 (2019): 86-100.
- [12] Alahmadi, Rehab, Chung Hyuk Park, and James Hahn. "Sequence-to-sequence image caption generator." In *Eleventh International Conference on Machine Vision (ICMV 2018)*, vol. 11041, p. 110410C. International Society for Optics and Photonics, 2019.
- [13] Hani, Ansar, Najiba Tagougui, and Monji Kherallah. "Image Caption Generation Using A Deep Architecture." In *2019 International Arab Conference on Information Technology (ACIT)*, pp. 246-251. IEEE, 2019.
- [14] Ballal, Noopur, and Sri Khetwat Saritha. "A Study of Deep Learning in Text Analytics." In *Social Networking and Computational Intelligence*, pp. 197- 205. Springer, Singapore, 2020.