

Image Captioning for Food Datasets: A Deep Learning Approach

Agustina Ghelfi

Universitat Autònoma de Barcelona
Barcelona, España

agustina.ghelfi@autonoma.cat

María José Millán

Universitat Autònoma de Barcelona
Barcelona, España

mariajose.millan@autonoma.cat

Laila Aborizka

Universitat Autònoma de Barcelona
Barcelona, España

lailamohamed.aborizka@autonoma.cat

Abstract

Image captioning bridges computer vision and NLP by generating textual descriptions from images. This project explores methods ranging from encoder-decoder models (ResNet + LSTM/GRU) to transformer-based architectures (ViT-GPT2, Qwen-2.5, ViT-Llama) using Food Ingredients and Recipes Dataset. We evaluate performance with BLEU, METEOR, and ROUGE metrics, highlighting the impact of model choice on caption quality. Our findings demonstrate the effectiveness of attention mechanisms and multimodal pretraining, while identifying challenges in generalization and evaluation.

Keywords: Image Captioning, Transformers, Multi-modal Models, Deep Learning.

1. Introduction

Image captioning is a fundamental task in computer vision and natural language processing that involves generating descriptive text based on visual input. It combines deep learning techniques from both domains, requiring models to extract meaningful visual features from images and translate them into coherent and contextually relevant text. Over the years, various approaches have been developed, ranging from traditional encoder-decoder architectures to more advanced methods incorporating attention mechanisms and transformer-based models.

The objective of this project is to explore different methodologies for image captioning, starting from the basics and gradually incorporating more sophisticated techniques. The project follows a structured learning path, beginning with data preprocessing and baseline models, and advancing towards modifications and improvements that align with state-of-the-art approaches. This structured ap-

proach allows us to assess model performance both quantitatively and qualitatively while iteratively improving the results.

For this study, we use the Food Ingredients and Recipes Dataset, which consists of 13,582 images of various dishes. The primary goal is to develop models capable of predicting dish titles based on image inputs. The dataset is divided into 80% training, 10% validation, and 10% test, ensuring a robust evaluation framework.

Throughout this project, we aim to gain deeper insights into the key components that contribute to effective image captioning systems and push the boundaries of traditional methods by exploring more advanced techniques.

2. Related work

Image captioning is a multidisciplinary research area that intersects computer vision and natural language processing (NLP) to generate textual descriptions of images. Over the years, various approaches have emerged, evolving from handcrafted rules and retrieval-based methods to modern deep learning-based architectures.

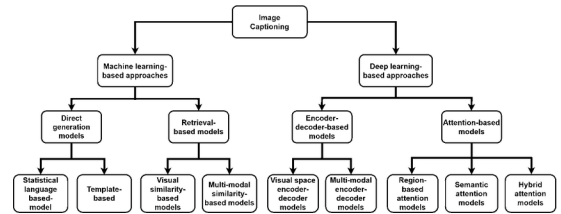


Figure 1. A comprehensive taxonomy of image captioning [7].

Early methods in image captioning primarily relied on template-based and retrieval-based techniques. In retrieval-based captioning, captions were either generated using pre-defined sentence structures or retrieved from an existing

pool based on image similarity. The process typically involved identifying visually similar images from a dataset, and using methods like Markov Random Fields and Lin similarity to measure semantic distance between the query and training images. Some approaches directly retrieved captions, while others generated new ones based on the retrieved images. Techniques like Analysis of Kernel Canonical Correlation and cosine similarity were used to rank and select the most relevant phrases [8]. However, these approaches lacked flexibility and struggled to generalize well to unseen images, often failing to produce meaningful captions for novel object combinations or complex scenarios. While they generated grammatically correct captions, they were limited in their ability to adapt to new image contexts, making them less effective for more diverse or dynamic image sets [10].

With the advent of deep learning, encoder-decoder architectures became the dominant paradigm. The first deep learning-based models utilized Convolutional Neural Networks (CNNs) for feature extraction and Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, for sequential text generation [11]. While these models significantly improved caption quality, they struggled with long-range dependencies and lacked a mechanism to focus on specific image regions.

To address this, attention mechanisms were introduced, allowing models to dynamically focus on relevant parts of an image while generating each word. Notable variants include soft attention, which learns a probabilistic weight distribution over image regions, and hard attention, which stochastically selects regions at each time step. More recently, self-attention and transformer-based architectures, such as Vision Transformers (ViTs) and multimodal transformers, have further improved captioning performance by capturing long-range dependencies and integrating vision and language more effectively [1, 10].

Another important development in image captioning is the use of reinforcement learning techniques, such as self-critical sequence training (SCST), which optimizes non-differentiable evaluation metrics like BLEU, CIDEr, and METEOR. These approaches enhance fluency and coherence beyond traditional maximum likelihood estimation [11]. Additionally, some works integrate graph-based representations and external knowledge bases to enrich captioning models with structured semantic relationships, leading to more contextually aware descriptions.

The field has also expanded into domain-specific applications, particularly in medical image captioning. Unlike general image captioning, where datasets such as MS COCO and Flickr30k are commonly used, medical image captioning requires specialized datasets like IU X-ray and MIMIC-CXR. These domain-specific models incorporate structured medical knowledge and leverage hierarchi-

cal learning techniques to generate accurate and clinically meaningful reports [1, 11].

Despite these advancements, several challenges remain. Current models still struggle with generating captions that are truly human-like in terms of coherence, diversity, and commonsense reasoning. Many image captioning models exhibit biases toward frequent objects and fail to describe rare or unseen entities effectively. Furthermore, evaluation metrics for image captioning remain a challenge, as automated scores often do not fully capture the quality of human-like descriptions [10].

Future research directions include enhancing multimodal pretraining techniques, leveraging large-scale vision-language models, and improving explainability in caption generation. The integration of external knowledge graphs, commonsense reasoning, and more sophisticated linguistic structures will also be crucial in advancing the field further.

3. Methodology

3.1. Base Architecture for Caption Generation

The baseline model was designed with an encoder-decoder architecture. In image captioning, the encoder is responsible for extracting meaningful visual representations from the input image, typically using CNNs to process and output high-level feature maps. The decoder interprets the visual features provided by the encoder and generates a corresponding natural language description, word by word. Recurrent neural networks (RNNs), are commonly employed for this task due to their ability to model sequential data. For the baseline model, the encoder utilized a ResNet-18, and the decoder was implemented with a GRU. Various alternatives for the encoder and decoder were explored to improve performance, and the results of these modifications were compared to find the most effective configuration.

3.1.1. Encoders

ResNet-18 Architecture. CNN composed of 18 layers, including convolutional layers and fully connected layers. It follows the original residual learning framework introduced by He et al., where identity shortcuts (skip connections) allow the model to learn residual functions instead of direct mappings [4]. This design significantly alleviates the vanishing gradient problem and enables the training of deeper networks. The architecture is structured as follows:

- An initial 7×7 convolutional layer with 64 filters and a stride of 2, followed by a 3×3 max pooling layer.
- Four residual stages, each consisting of 2 basic residual blocks. Each basic block includes two 3×3 convolutional layers with batch normalization and ReLU activation.
- A global average pooling layer at the end aggregates spatial information.
- A fully connected layer produces the final classification scores.

ResNet-50 Architecture. Deeper and more expressive variant than ResNet-18, comprising 50 layers. It utilizes bottleneck residual blocks to enhance representational capacity while maintaining computational efficiency [4]. Its design reduces the computational cost while enabling deeper models:

- Initial 7×7 convolution and max pooling.
- Four stages, each containing a different number of bottleneck blocks, with the final stage having 2048 output filters. Each bottleneck block has:
 - A 1×1 convolution to reduce dimensionality.
 - A 3×3 convolution to process spatial features.
 - Another 1×1 convolution to restore dimensionality.

3.1.2. Decoder

LSTM Architecture. The Long Short-Term Memory network is a type of recurrent neural network designed to better capture long-term dependencies by mitigating the vanishing gradient problem [2]. It does so through a memory cell and a system of gating units that regulate the flow of information. The **input gate** controls how much of the new input is added to the cell state, the **forget gate** determines the amount of past information to retain or discard from the cell state, and the **output gate** regulates the exposure of the internal cell state to the output.

The decoder is composed of one or more stacked LSTM layers, each processing the output of the previous layer. During decoding, the final hidden state from the encoder is used to initialize the LSTM decoder's hidden and cell states. The output at each time step is passed through a linear transformation followed by a softmax function to produce a probability distribution over the target vocabulary.

GRU Architecture. The Gated Recurrent Unit is a variant of the LSTM that combines the memory and gating mechanisms into a simpler structure [2]. It was introduced to improve training efficiency while maintaining competitive performance. Each GRU unit contains two gates, an **update gate** that controls the degree to which the unit updates its activation or content, and a **reset gate** that determines how much of the past information to forget.

The decoder uses stacked GRU layers, with each layer processing sequential information and producing outputs passed to the next layer. Initial hidden states for the decoder are typically derived from the final encoder states. The output of each time step is projected through a linear layer followed by softmax activation for token prediction.

Transformer Decoder Architecture. Provides a non-recurrent alternative capable of modeling long-range dependencies through attention mechanisms. Its architecture consists of several key components:

- **Self-attention mechanism:** allows the decoder to weigh the importance of previous tokens in the sequence, ensuring that each token is generated in the context of the entire caption generated so far.
- **Decoder layers:** Each layer of the Transformer Decoder includes multi-head self-attention, a feed-forward network, residual connections, and layer normalization. These components help capture long-range dependencies in the generated sequence.
- **Positional encoding:** Since the Transformer lacks a built-in notion of token order, positional encodings are added to the input embeddings to inform the model of the token's position in the sequence.

Teacher Forcing: A commonly used technique during the training of decoders based on RNNs is *teacher forcing*. In this approach, during training, instead of using the model's previous predictions to generate the next word, the actual word from the training set is fed as input for the next prediction. This helps the model by providing more accurate information, preventing the accumulation of errors over time. Also allows for faster convergence and improves performance during training. However, it can introduce a discrepancy between training and inference behavior, as during inference, the model must rely solely on its own predictions.

3.2. ViT-GPT2 Model

Vision Transformer (ViT) serves as the visual encoder, extracting image features, while GPT-2 acts as the text decoder, generating captions in natural language.

ViT Architecture: Model for image processing based on the Transformer architecture, transforms an image into a sequence of patches, processes it using a Transformer to capture global relations, and uses a classification token for image representation [3].

- **Image preprocessing:** divides the image into non-overlapping patches, flattens them, and projects them into a fixed-dimensional embedding space. A special classification token precedes the sequence to serve as a global representation of the image. Positional embeddings are also added to retain spatial information across patches.
- **Transformer encoder:** Composed of several layers that include multi-head self-attention and feedforward MLP blocks. Each layer applies Layer Normalization before the main operation and uses residual connections afterward. The final output of the classification token serves as the overall image representation.
- **Inductive Bias:** incorporates weaker inductive biases compared to CNNs. It does not inherently capture local structures or enforce translation invariance, instead relying on the patch-based input and positional embeddings to model spatial relationships.

GPT-2 Architecture: Decoder-only Transformer model for text generation, trained to predict the next word in a sequence using causal self-attention and MLP layers for processing concept vectors [5].

- **Word Embeddings:** Words are mapped to high-dimensional vectors that are nearly orthogonal to each other, acting as concept vectors.
- **Causal Self-Attention:** Each position can only attend to previous ones, ensuring that predictions are not influenced by future tokens.
- **MLP Layers:** Projects the embeddings to a higher dimension and back. This step captures non-linear relationships between concepts.
- **Word Prediction:** The final output is projected back to the word embeddings space, and a dot product is computed with each word embedding to predict the next token using a softmax function.

3.3. Qwen-2.5 7B Multimodal

QWEN 2.5 is a Transformer-based model optimized for complex NLP and multimodal tasks.

- **Self-Attention:** Multi-head self-attention captures long-range dependencies, with dynamic scaling for improved focus on relevant tokens.
- **Adaptive Attention Span:** Dynamically adjusts attention span for efficiency, reducing redundant computations.
- **Position Encoding:** Combines absolute/relative position encodings and Rotary Position Embeddings (RoPE) for better sequence handling.
- **Multi-Modal Processing:** A cross-attention mechanism integrates text and images using dual encoders (ViT for images, Text Transformer for language).
- **Efficient Training & Inference:** Sparse attention and gradient checkpointing reduce computational cost, with Flash Attention for faster inference on accelerators.

3.4. Vit-Llama 3.2 Model

The Llama 3.2 model is a state-of-the-art language model developed by Meta, designed for high-performance, large-scale language understanding and generation tasks. It is built on a transformer-based architecture similar to GPT-3, optimized for efficiency and scalability, and features multiple layers of self-attention to capture long-range dependencies and contextual relationships. This model can be fine-tuned to generate captions by conditioning on visual embeddings from the ViT encoder. It follows a cross-modal learning approach, where visual features guide the language generation process. Llama 3.2 integrates with the ViT encoder, allowing the model to handle multimodal inputs, using visual features as input to generate accurate captions aligned with the visual content. In this project, both the 1B and 3B parameter versions of Llama 3.2 were explored.

4. Experimental Settings

4.1. Dataset

The **Food Ingredients and Recipes Dataset** (licensed under CC BY-SA 3.0) contains 13,582 images of food dishes [9]. Each image has a corresponding entry in a CSV file, which has the following columns:

- **Title:** The dish title.
- **Ingredients:** The list of ingredients (raw data).
- **Instructions:** Cooking instructions for the dish.
- **Image Name:** Corresponding image filename.
- **Cleaned Ingredients:** Processed ingredients list.

Data cleaning involved removing 5 rows with missing titles and deleting rows with invalid image names, resulting in 13,501 valid image-caption pairs. Figure 2 shows a word cloud analysis of the dish titles, which reveals frequent terms such as “Salad” and “Chicken”.



Figure 2. Word cloud analysis of the ‘Food Ingredients and Recipes Dataset’.

4.2. Metrics

- **BLEU** (Bilingual Evaluation Understudy) is a metric for evaluating machine-generated translations based on n-gram precision, with a modified unigram precision. It includes a brevity penalty to balance word choice, order, and length, penalizing overly long translations while rewarding more balanced outputs. For this project, we used BLEU-1 and BLEU-2, which measure the precision of unigrams (single words) and bigrams (two consecutive words), respectively, between the generated and reference captions.
- **METEOR** (Metric for Evaluation of Translation with Explicit ORdering): enhances BLEU by emphasizing recall, considering exact matches, stemmed words, and synonyms between candidate and reference sentences. It is designed to provide a more holistic evaluation by capturing semantic equivalences and not just n-gram matches [6].
- **ROUGE** (Recall-Oriented Understudy for Gisting Evaluation): initially developed for text summarization, focuses on n-gram recall, measuring the overlap between candidate and reference n-grams. It is particularly useful for tasks where recall of important words and phrases is crucial [6].

4.3. Implementation Details

4.3.1. Environment and Setup

The programming language used for this project is Python (version 3.13.2), with PyTorch (version 2.6.0+cu124) and Hugging Face Transformers (version 4.49.0) as the primary frameworks. The hardware consists of an NVIDIA RTX 3090 GPU, with 24 GB of GPU memory. The operating system is Ubuntu 20.04.6 LTS (Focal Fossa).

4.3.2. Model-Specific Details

A. Base Architecture for Caption Generation The encoder can be configured to use either ResNet-18 or ResNet-50. The decoder can be configured as either GRU, LSTM or Transformer Decoder. During training, teacher forcing can be applied, where the actual ground-truth token is fed into the decoder at each time step instead of its previous prediction. This helps improve convergence and reduces error accumulation.

Depending on configuration, the optimizer can be Adam, AdamW, or SGD, with hyperparameters (learning rate, weight decay) specified externally. An early stopping mechanism monitors the validation loss (with a delta threshold of 0.001) to prevent overfitting.

B. ViT-GPT2 Model Images are preprocessed using a `ViTImageProcessor`, which resizes, normalizes, and splits them into patches before passing them through the encoder to extract features. The text decoder is integrated via a `VisionEncoderDecoderModel`, and captions are tokenized using the corresponding `AutoTokenizer`. GPT-2 generates text autoregressively, enabling caption generation conditioned on visual features.

The model supports selective freezing to reduce trainable parameters. In the ViT encoder, patch embeddings and any number of initial layers can be frozen. Likewise, in the GPT-2 decoder, token (wte) and positional (wpe) embeddings, as well as early transformer blocks, can be optionally frozen. This flexibility helps focus learning on task-specific features.

Training uses the AdamW optimizer with a learning rate of $2e-5$. Batch sizes and number of epochs are set via command-line arguments. Early stopping (patience of 5 epochs, threshold of 0.001) is applied to prevent overfitting. A `CustomTrainer` handles training, extending the Hugging Face `Trainer` to compute loss, return model outputs if needed, and apply different learning rates to encoder and decoder.

C. Qwen-2.5 7B Multimodal For caption generation, the `Qwen2_5_VLForConditionalGeneration` class is used alongside its corresponding `AutoProcessor`, which applies a chat template and handles image inputs

using custom preprocessing functions. Images are loaded from a predefined dataset and formatted into message structures that include both the image and a short textual prompt.

During inference, messages are tokenized, images are preprocessed, and captions are generated autoregressively with a maximum of 128 tokens. The model outputs a single descriptive phrase per image, which is then cleaned and stored in a CSV file. This setup allows the model to generate visually grounded text without requiring additional fine-tuning, leveraging the instruction-following capabilities of the base model.

D. ViT-Llama 3.2 Model The ViT encoder processes images by dividing them into 16×16 patches and extracting visual features using a 12-layer transformer with a hidden size of 768 and 12 attention heads. For text generation, Llama 3.2 is used in its 1B and 3B parameter variants, both configured as causal language models with 12 layers and a hidden size of 1024.

To integrate both modalities, the visual features extracted by ViT are projected into the Llama embedding space and concatenated with text embeddings, enabling joint processing by the decoder. The model is fine-tuned using LoRA (Low-Rank Adaptation), which introduces low-rank updates to reduce the number of trainable parameters. Specifically, LoRA is applied with rank 4, alpha 16, and dropout 0.1 for the 1B version, and with rank 4, alpha 8, and dropout 0.1 for the 3B version. Training is performed with the AdamW optimizer at a learning rate of $1e-4$, over 10 epochs, with a batch size of 4 for training and 2 for validation. The combined ViT and Llama components are implemented in the custom `ViTLlamaForCaptioning` class, designed to handle multimodal input efficiently.

References

- [1] Lakshita Agarwal and Bindu Verma. From methods to datasets: A survey on image-caption generators. *Multimedia Tools and Applications*, 83:28077–28123, 2024. 2
- [2] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014. 3
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. arXiv:2010.11929. 3
- [4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 770–778, 2016. 2, 3

- [5] Johannes Knittel, Tushaar Gangavarapu, Hendrik Strobelt, and Hanspeter Pfister. Gpt-2 through the lens of vector symbolic architectures. In *2nd Workshop on Attributing Model Behavior at Scale (ATTRIB), NeurIPS*, 2024. Workshop paper. 4
- [6] Sara Sarto, Marcella Cornia, and Rita Cucchiara. Image captioning evaluation in the age of multimodal llms: Challenges and future perspectives, 2025. 4
- [7] Himanshu Sharma and Devanand Padha. A comprehensive survey on image captioning: from handcrafted to deep learning-based techniques, a taxonomy and open research issues. *Artificial Intelligence Review*, 56:1–43, 2023. 1
- [8] Himanshu Sharma, Manmohan Agrahari, Sujeet Kumar Singh, Mohd Firoj, and Ravi Kumar Mishra. Image captioning: a comprehensive survey. In *2020 International Conference on Power Electronics & IoT Applications in Renewable Energy and its Control (PARC)*, pages 325–328. IEEE, 2020. 2
- [9] Manan Sharma. Food ingredients and recipe dataset with images. <https://www.kaggle.com/datasets/pes12017000148/food-ingredients-and-recipe-dataset-with-images>, 2020. Accessed: 2025-04-05. 4
- [10] Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, Silvia Cascianelli, Giuseppe Fiameni, and Rita Cucchiara. From show to tell: A survey on deep learning-based image captioning. *arXiv preprint*, 2021. 2
- [11] Liming Xu, Quan Tang, Jiancheng Lv, Bochuan Zheng, Xi-anhua Zeng, and Weisheng Li. Deep image captioning: A review of methods, trends, and future challenges. *Neuro-computing*, 546:126287, 2023. 2