

Article

Making Images Speak: Human-Inspired Image Description Generation

Chifaa Sebbane *¹, Ikram Belhajem ¹ and Mohammed Rziza *

LRIT, Faculty of Sciences in Rabat, Mohammed V University in Rabat, Rabat 10000, Morocco

* Correspondence: chifaa_sebbane@um5.ac.ma (C.S.); m.rziza@um5r.ac.ma (M.R.)

Abstract: Despite significant advances in deep learning-based image captioning, many state-of-the-art models still struggle to balance visual grounding (i.e., accurate object and scene descriptions) with linguistic coherence (i.e., grammatical fluency and appropriate use of non-visual tokens such as articles and prepositions). To address these limitations, we propose a hybrid image captioning framework that integrates handcrafted and deep visual features. Specifically, we combine local descriptors—Scale-Invariant Feature Transform (SIFT) and Bag of Features (BoF)—with high-level semantic features extracted using ResNet50. This dual representation captures both fine-grained spatial details and contextual semantics. The decoder employs Bahdanau attention refined with an Attention-on-Attention (AoA) mechanism to optimize visual-textual alignment, while GloVe embeddings and a GRU-based sequence model ensure fluent language generation. The proposed system is trained on 200,000 image-caption pairs from the MS COCO train2014 dataset and evaluated on 50,000 held-out MS COCO pairs plus the Flickr8K benchmark. Our model achieves a CIDEr score of 128.3 and a SPICE score of 29.24, reflecting clear improvements over baselines in both semantic precision—particularly for spatial relationships—and grammatical fluency. These results validate that combining classical computer vision techniques with modern attention mechanisms yields more interpretable and linguistically precise captions, addressing key limitations in neural caption generation.

Keywords: image captioning; deep learning; visual attention; neural networks; ResNet; GloVe embeddings; SIFT; Bag of Features; CIDEr; assistive technologies; MS COCO; NLP



Academic Editors: Wei Zhou, Guanghui Yue, Wenhan Yang and Gholamreza Anbarjafari (Shahab)

Received: 21 December 2024

Revised: 26 February 2025

Accepted: 20 March 2025

Published: 28 April 2025

Citation: Sebbane, C.; Belhajem, I.; Rziza, M. Making Images Speak: Human-Inspired Image Description Generation. *Information* **2025**, *16*, 356. <https://doi.org/10.3390/info16050356>

Copyright: © 2025 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Automatically generating natural language descriptions from visual content, known as image captioning (IC), is a fundamental challenge at the intersection of computer vision and natural language processing. As the volume of visual content continues to expand across diverse domains, IC plays an increasingly vital role in enabling machines to **understand and communicate** visual information in a way that is easily understandable by humans. Its applications span a wide range of fields, including accessibility for the visually impaired. This is particularly relevant when dealing with images captured by blind users in **unconstrained, real-world settings** [1]. Other important applications include autonomous driving [2], content-based image retrieval [3], and interaction between humans and artificial intelligence.

Early approaches to image captioning (IC) relied on handcrafted visual features and statistical models. Signal processing techniques, such as Fourier transforms [4] and wavelets [5], were used to encode image content at low resolution and limited detail, while Hidden Markov Models (HMMs) [6] was employed to model language.

Descriptors like SIFT [7] and Bag of Features (BoF) [8] contributed to improved object and texture recognition. However, these methods often struggled to capture complex semantic relationships and lacked generalization capabilities in diverse or cluttered scenes [9]. These limitations underscored the need for more adaptive, data-driven approaches, paving the way for neural architectures.

Hybrid and early neural models sought to bridge vision and language more effectively. BabyTalk [10] combined symbolic visual inputs with structured sentence templates, while the Log-Bilinear Model (LBL) [11] and its multimodal extensions [12,13] began integrating learned embeddings across modalities. These models improved fluency and multimodal representation but lacked strong alignment mechanisms and struggled with visual–textual complexity.

The deep learning era brought major advances. Convolutional Neural Networks (CNNs) [14] became standard for visual encoding, while RNNs, notably LSTM [15] and GRU [16], enabled sequential caption generation [17]. Models like m-RNN [18] and “Show and Tell” [19] established strong baselines, but still had difficulty modeling long-range dependencies and complex scene dynamics.

To improve specificity, object detection was integrated into captioning systems. Approaches such as R-CNN [20], Fast R-CNN [21], Faster R-CNN [22], and YOLO [23] enabled models to identify and focus on salient visual regions. Karpathy and Fei-Fei [24] further enhanced semantic grounding by aligning textual fragments with specific image regions using multimodal embeddings. Yet these methods heavily depend on accurate detection and predefined bounding boxes, which can limit flexibility in dynamic or unconstrained environments.

Attention mechanisms addressed these issues by allowing models to dynamically attend to relevant regions during caption generation [25]. Extensions such as AoA [26] improved spatial alignment. However, these techniques may still overlook subtle or non-salient elements, especially when visual features lack granularity.

More recently, hybrid architectures combining object-level attention and CNN encoding have been proposed [27], though they often neglect handcrafted descriptors and remain dependent on object detection accuracy. Other alternatives include weakly supervised techniques such as Multiple Instance Learning (MIL) [28], and structured representations via scene graphs [29,30] or image segmentation [31,32]. Scene graphs explicitly capture object relationships and spatial arrangements, enabling more structured and relational caption generation. Likewise, segmentation approaches help delineate coherent visual regions, either through semantic labels or unsupervised grouping, thereby improving visual grounding and supporting the generation of more semantically aligned descriptions. However, these methods often rely on complex annotations or demand substantial preprocessing.

Transformer-based models like ViT [33] and CPTR [34] perform well by leveraging global visual context. Meanwhile, VLP methods such as CLIP [35] and GPT-4V [36] show promise for general-purpose multimodal reasoning. Even so, these systems typically require vast data and computational resources, and often struggle to generalize to out-of-distribution (OOD) samples [37].

As a more flexible and lightweight alternative, recent architectures combine a CNN encoder (e.g., ResNet) with an LSTM decoder and integrate a pretrained auxiliary language model such as BERT. The initial captions are typically generated by a baseline visual captioning model such as Show, Attend, and Tell, based on visual features extracted by the CNN and decoded sequentially by the LSTM. These captions are then refined through a dedicated fusion module that merges contextual representations from the language model with the visual decoding stream [38]. Using strategies like Deep Fusion, Cold Fusion, or

Hierarchical Fusion, this architecture effectively corrects syntactic and semantic errors while requiring less data and computational cost. Moreover, it provides a flexible framework adaptable to domain-specific applications such as medical or assistive image description.

Several key challenges therefore remain open: (1) distinguishing between visually grounded and contextually inferred words; (2) preserving fine-grained spatial information; and (3) reducing the computational cost of Transformer-based architectures.

In this paper, we propose a hybrid image captioning model that combines handcrafted local descriptors, such as SIFT and Bag of Features (BoF), with global deep features extracted from a ResNet50 backbone. This dual representation captures both low-level structural detail and high-level semantic context. The decoder integrates a dual-attention mechanism, combining Bahdanau attention and Attention on Attention (AoA), while a GRU-based language model, enhanced with pre-trained GloVe embeddings [39], ensures fluent and coherent caption generation. Compared to standard CNN-based architectures, our model enhances interpretability and linguistic expressiveness, all while maintaining a low computational footprint.

The system is developed and evaluated on extensive benchmarks, including 200,000 image–caption pairs from the MS COCO train2014 split, 50,000 held-out validation samples, and the Flickr8K dataset. The model demonstrates competitive performance, with a CIDEr score of 128.3, validating the benefits of our hybrid architecture in capturing both spatial detail and linguistic nuance.

The remainder of this paper is structured as follows: Section 2 presents the proposed methodology. Section 3 discusses the experimental results. Finally, Section 4 concludes the paper and outlines future research directions.

2. Materials and Methods

Conventional CNN-based feature extraction methods often rely on large annotated datasets, suffer from loss of fine spatial details, and exhibit limited generalization to unseen objects. To address these limitations, this research proposes a hybrid model that combines classical and deep learning techniques for image captioning.

2.1. Image Feature Extraction

Our image captioning pipeline draws on both classical and deep learning features to capture fine-grained details as well as broad semantic cues. By integrating a Bag of Features (BoF) representation (derived from SIFT) with ResNet50’s residual architecture [40], we obtained a unified feature set capturing both low-level textures and high-level semantic cues. This hybrid approach ensures that both local structural information (through SIFT) and global contextual understanding (through ResNet50) contribute to the captioning process.

2.1.1. SIFT-Based Local Feature Extraction and Bag of Features Representation

SIFT for Local Feature Extraction

SIFT was employed to detect robust and spatially precise keypoints, preserving the fine-grained structural details crucial for attention mechanisms. In our implementation, keypoints were refined using contrast and curvature thresholding to discard low-contrast and edge-sensitive candidates, thereby retaining only the most salient spatial information for effective object localization.

Descriptor Computation and BoF

To convert the resulting variable-length SIFT descriptors into a fixed-length vector, we clustered them into “visual words” using k -Means [41]. This unsupervised approach

eliminates the need for fine-tuning or annotated labels, and improves generalization by grouping common visual motifs even for previously unseen objects while also reducing dataset bias. Although k -Means clustering is computationally efficient with a complexity of $\mathcal{O}(nkd)$ (where n is the number of descriptors, k the number of clusters, and d the descriptor dimension), its sensitivity to initialization can sometimes lead to suboptimal clustering. To mitigate this, we used a balanced initialization strategy to improve robustness. By constructing a histogram of visual word frequencies (length $k = 100$), we obtained BoF vectors that effectively capture localized shape and texture information essential for captioning tasks.

2.1.2. ResNet50-Based Global Feature Extraction

To capture high-level semantic information, we used ResNet50, pre-trained on ImageNet. Initially, we experimented with VGG16 [42], but ResNet50 demonstrated superior performance due to its residual connections, which mitigate vanishing gradients and allow deeper feature extraction. We did not evaluate EfficientNet or Transformers due to computational constraints and because ResNet50 provided an optimal trade-off between model complexity and feature richness for our dataset.

Feature Processing Pipeline

- **Preprocessing:** All images are resized to 256×256 and normalized according to ResNet50's training specifications.
- **Extraction:** Using 'include_top=False', we discard the fully connected layers and retain the final convolutional block, which produces an $(8 \times 8 \times 2048)$ feature map encapsulating global semantics.

2.1.3. Feature Fusion and Alignment for Image Captioning

To generate comprehensive image captions, we fused the local (SIFT + BoF) and global (ResNet50) features through a *tiling-and-concatenation* approach:

1. **Tiling (Alignment):** The 1D BoF histogram is replicated to match ResNet50's spatial resolution (e.g., 8×8), preserving its original frequency distribution without introducing additional learnable parameters.
2. **Channel Concatenation:** The tiled BoF is concatenated along the channel dimension of the ResNet50 feature map (e.g., $2048 + 100 = 2148$ channels), merging local descriptors and global context into a single 3D tensor.
3. **Reshaping for Decoder Input:** The resulting $(8 \times 8 \times 2148)$ tensor is then flattened into a (64×2148) matrix to ensure compatibility with the attention-based decoder.

This fusion method is visualized in Figure 1.

2.2. Data Preparation and Batch Construction

To enable efficient training of the image captioning model, we developed a streamlined data pipeline that aligns image features with their corresponding captions, transforms them into a numerical format, and constructs batches optimized for GPU processing. This section details each stage of the data flow.

2.2.1. Image–Caption Mapping

Each image in the MS COCO train2014 dataset is associated with up to five captions describing its content. Using the official annotation file, we parsed the JSON content and built a dictionary (`image_id_to_captions`) that links each image ID to its corresponding raw captions. This many-to-one mapping allows the model to learn diverse linguistic expressions for the same visual content.

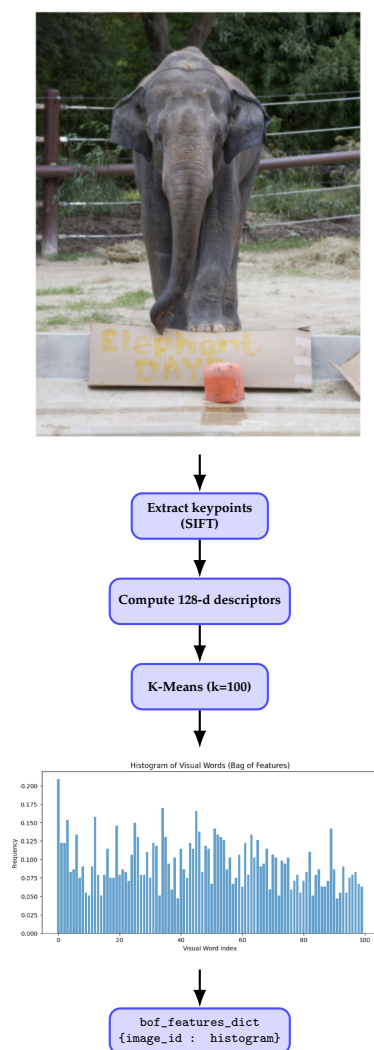


Figure 1. SIFT + BoF Feature Extraction Pipeline: (1) SIFT detects local keypoints and computes 128-dimensional descriptors. (2) K-Means clusters these descriptors into visual words. (3) A histogram of visual word frequencies is generated. (4) The resulting histogram is stored in a dictionary named `bof_features_dict`, indexed by each image’s `image_id`, for subsequent retrieval.

The MS COCO train2014 set was used for both training and validation via an internal split (80–20). The Flickr8K dataset, entirely unseen during training, was reserved solely for inference-based generalization testing (see Section 3).

2.2.2. Text Preprocessing and Tokenization

We utilized the Keras Tokenizer from `tensorflow.keras.preprocessing.text` to transform raw captions into structured numerical inputs.

The Tokenizer played a dual role in our pipeline:

- **Before training:** It converted raw captions into structured numerical sequences, ensuring consistent word representation.
- **During caption generation:** When the decoder predicted a word index, the Tokenizer mapped it back to the corresponding word to reconstruct a human-readable sentence.

The complete pipeline, including all key preprocessing steps applied to the captions, is visualized in Figure 2. While operations such as lowercasing and tokenization may reduce linguistic variety, we counter this effect by assigning multiple captions to each image. This strategy introduces diverse expressions for the same content, enabling the model to learn richer and more flexible language patterns.

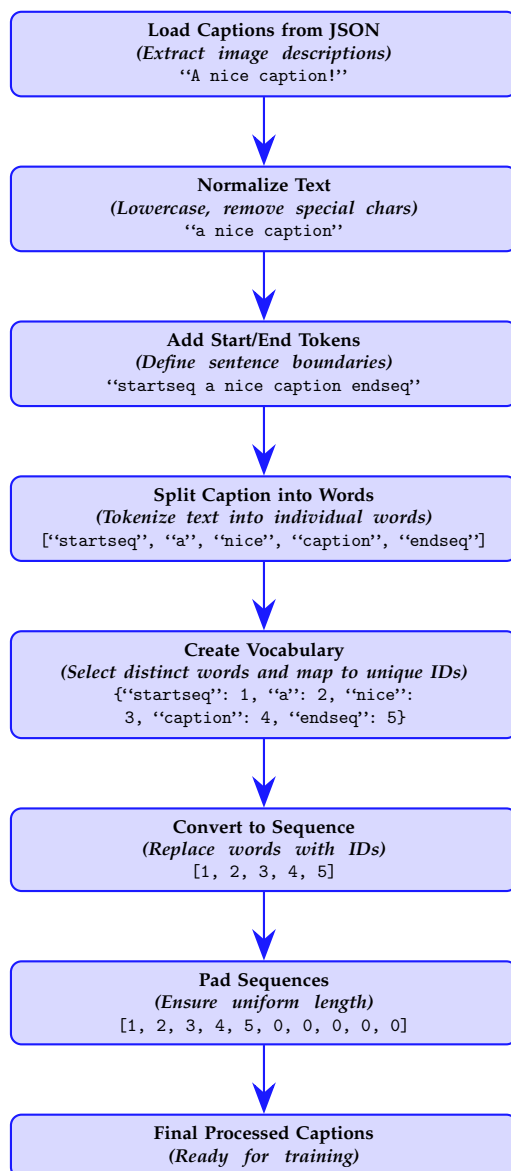


Figure 2. Text preprocessing pipeline for caption modeling. Steps include normalization, tokenization, vocabulary indexing, and padding.

2.2.3. Feature–Caption Alignment

Each image ID was linked to a precomputed image feature vector from the hybrid SIFT + ResNet50 model (`feature_dict`). These vectors were paired with each of their corresponding tokenized captions to ensure proper alignment between visual and linguistic inputs during supervised learning.

2.2.4. Train–Test Splitting

We performed an 80–20 split of the MS COCO `train2014` set:

- **Training set:** 200,000 image–caption pairs.
- **Validation set:** 50,000 image–caption pairs.

The Flickr8K dataset was used only for generalization evaluation on unseen images.

2.2.5. Optimized Batch Pipeline with TensorFlow

To support efficient GPU-based training, we implemented a high-throughput batch construction system using TensorFlow:

- **Mapping:** A custom function, `map`, is used to retrieve image features from the `feature_dict` and corresponding captions, casting all tensors to `float32` for GPU compatibility.
- **Dataset construction:** A generator pipeline is built using `tf.numpy_function`, allowing for efficient data processing and integration with TensorFlow's data pipeline.
- **Shuffling and batching:** The data is shuffled with a buffer size of 1000 and batched with a size of 64 to ensure balanced and efficient training.
- **Prefetching:** `tf.data.experimental.AUTOTUNE` is used to minimize idle GPU time and ensure that data is prepared ahead of time for the training process.

Batch Shapes

- **Image features:** $(64, 64, 2148)$. Each image is divided into 64 spatial patches (e.g., following an 8×8 grid), and a hybrid feature vector of 2148 dimensions is extracted for each patch. This high-dimensional representation results from the concatenation of handcrafted local descriptors (SIFT) with deep visual features obtained from a pretrained ResNet50. By combining both local texture cues and semantic content, this representation provides a rich multimodal embedding of the visual input. Hence, the resulting shape per batch is $(64, 64, 2148)$, where 64 is the batch size, 64 the number of patches per image, and 2148 the feature dimensionality per patch.
- **Captions:** $(64, 51)$, representing the tokenized and padded textual sequences corresponding to each image.

This pipeline ensures synchronized delivery of visual and textual data, while enabling smooth, scalable training across hardware resources.

2.3. Model Details

Figure 3 illustrates the end-to-end structure of the image captioning system, showcasing both the encoder–decoder pipeline and the integration of attention mechanisms.

2.3.1. Encoder Architecture and Functionality

The encoder architecture, illustrated in Figure 4, is implemented using TensorFlow Keras's subclassing API by extending the `Model` base class. This approach offers flexibility and facilitates integration within the TensorFlow ecosystem. The encoder is designed to merge deep and handcrafted visual features and transform them into a compact and expressive representation for subsequent processing by the decoder.

Dense Layer Functionality

As shown in Figure 4, the encoder takes as input a feature map of shape $(batch_size, 64, 2148)$, resulting from the concatenation of ResNet50 features and SIFT+BoF descriptors. This is projected into a lower-dimensional space (`embed_dim = 200`) using a Dense layer with Leaky ReLU activation. The output $(batch_size, 64, 200)$ retains key spatial and semantic information in a more compact form.

Leaky ReLU Activation

To mitigate the “dying ReLU” problem—where neurons become inactive and output zero for all inputs due to zero gradients in the negative domain—a Leaky ReLU activation function is used. This variant introduces a small non-zero gradient (`alpha = 0.01`) for negative inputs, ensuring continuous weight updates and improved model stability during training [43].

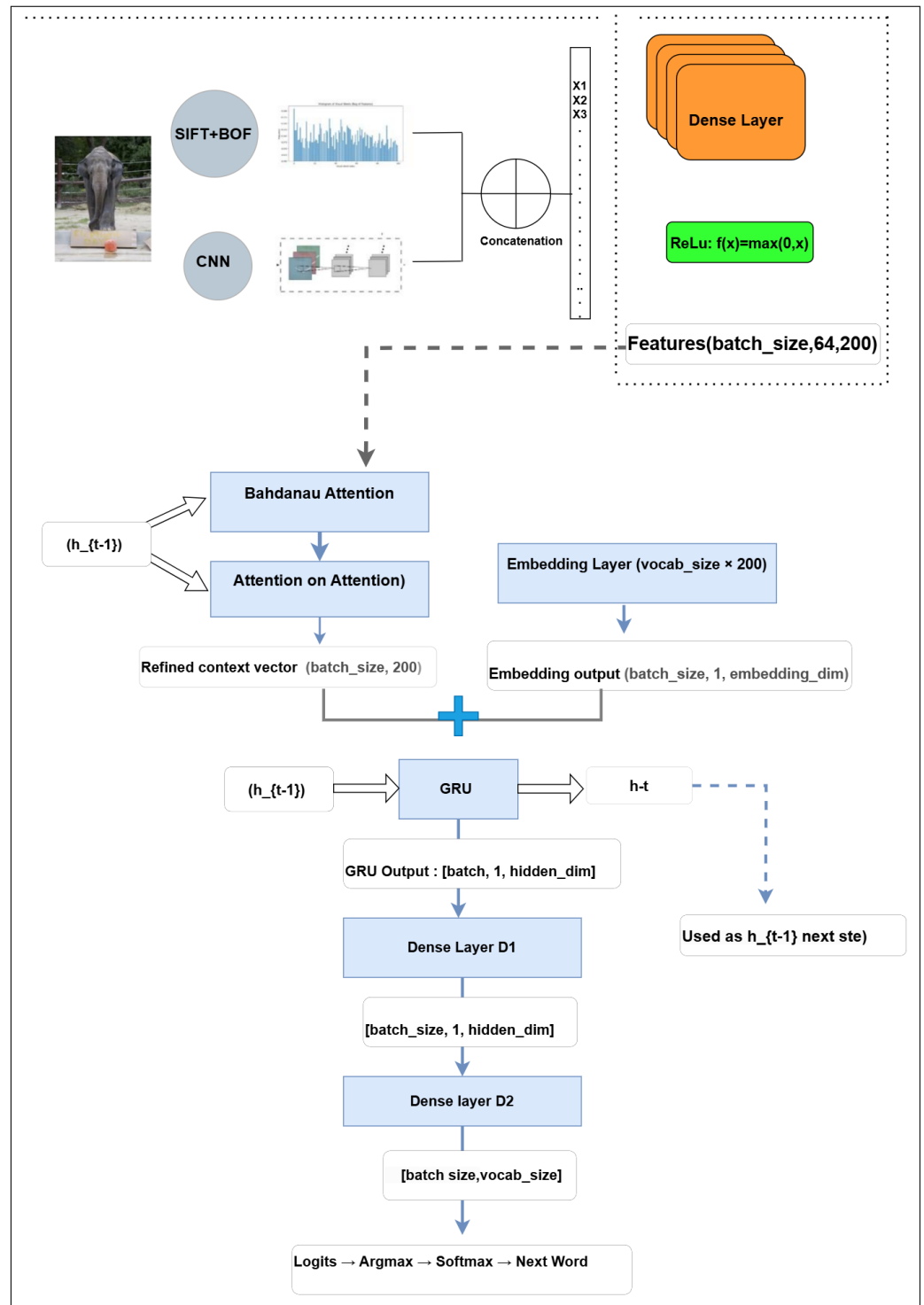


Figure 3. Architecture of the proposed image captioning model. The model integrates handcrafted features (SIFT+BoF) with deep visual features (ResNet50). These are concatenated and passed through a dense layer with Leaky ReLU activation. The decoder incorporates Bahdanau and AoA (Attention on Attention) mechanisms. Word embeddings are initialized using GloVe vectors, and a GRU-based decoder generates captions, followed by two Dense layers (D1 and D2) for final prediction.

Role in Caption Generation

The encoder transforms the fused visual inputs into lower-dimensional feature sequences that are structured for attention-based decoding. Its role is to condense both local and global cues into a compact format that enables precise focus during caption generation, while maintaining the spatial layout necessary for alignment with textual tokens.

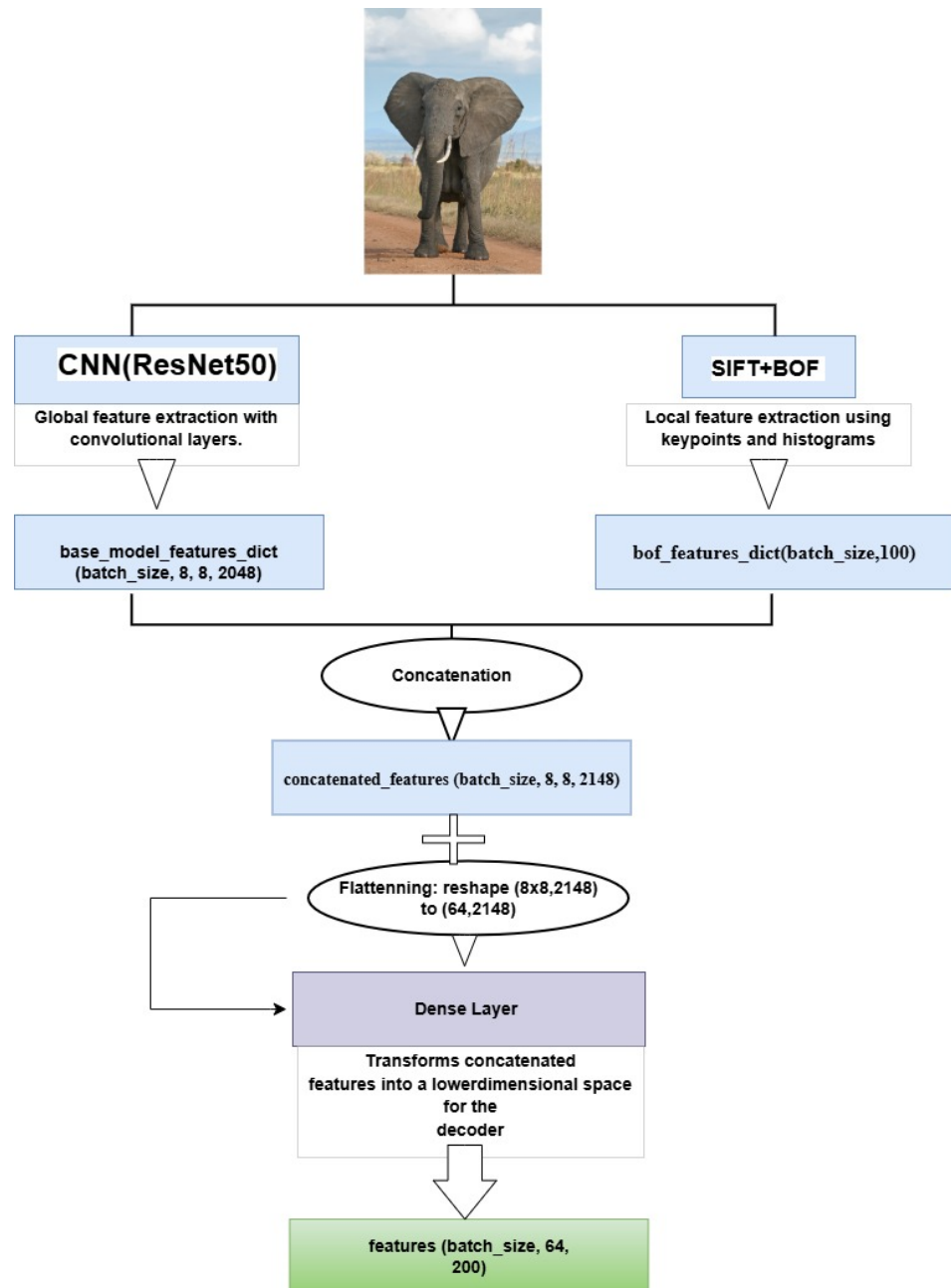


Figure 4. Encoder architecture for image feature processing. The encoder fuses global visual features extracted from ResNet50 with local descriptors obtained via SIFT and Bag-of-Features (BoF). The combined feature representation is passed through a dense layer to produce compact embeddings suitable for decoding.

2.3.2. Decoder Architecture

The decoder transforms the encoded visual features into textual descriptions by combining attention mechanisms with sequential modeling. This design enhances both contextual relevance and fluency during caption generation.

Figure 5 illustrates the internal components of the decoder, which combines dual attention strategies with GRU-based sequence generation.

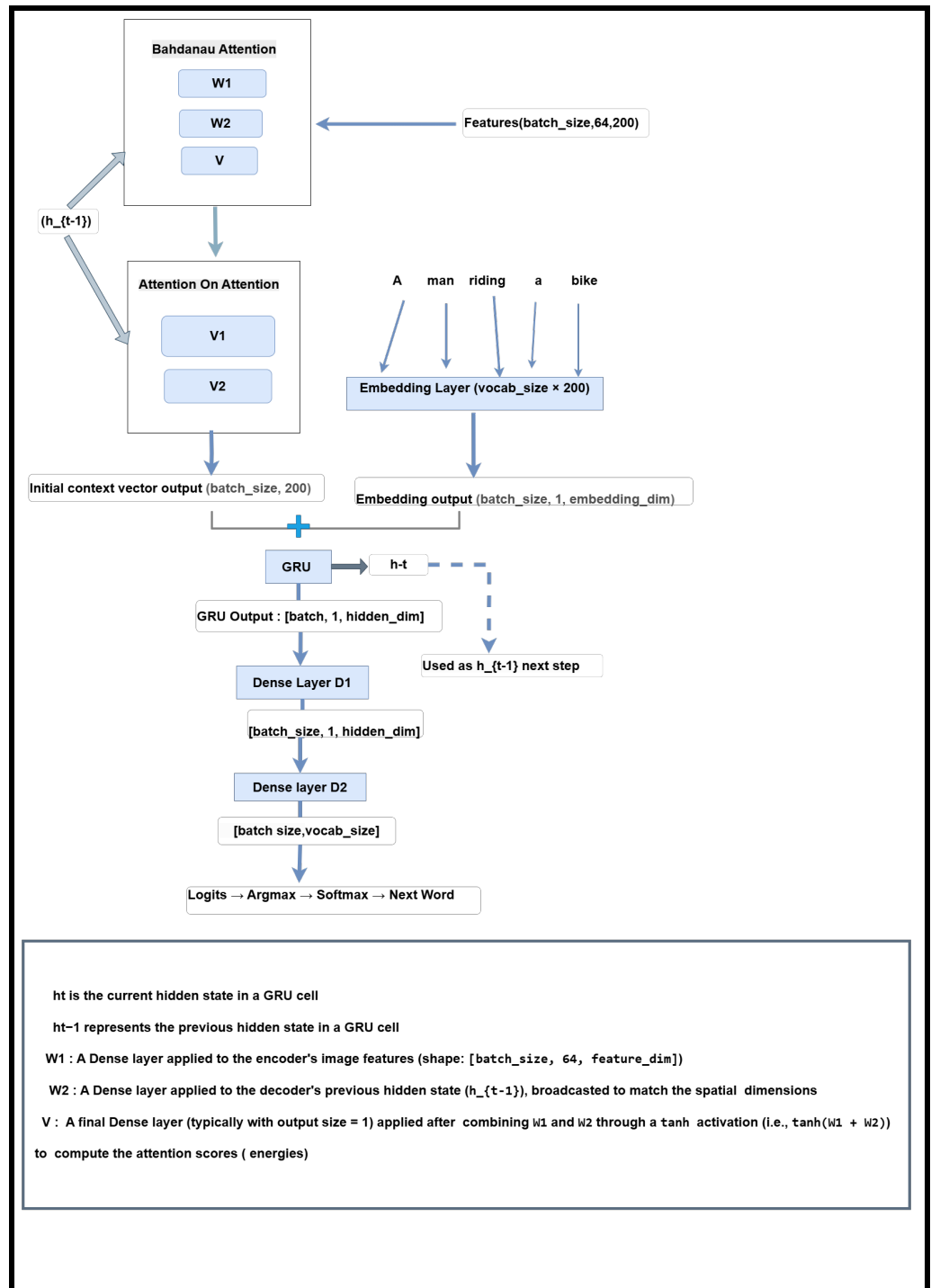


Figure 5. Decoder architecture for image caption generation. The decoder integrates a pre-trained embedding layer, Bahdanau attention followed by Attention-on-Attention (AoA), a GRU unit, and two dense layers. At each decoding step, the context vector is computed by applying Bahdanau attention over the image features using the decoder’s previous hidden state h_{t-1} . This context vector is then refined by the AoA module and concatenated with the embedding of the current input word. The resulting representation is passed to the GRU, which outputs an updated hidden state h_t . This hidden state is subsequently processed by two dense layers to produce a probability distribution over the vocabulary, enabling the prediction of the next word.

Embedding Layer Preparation Using Pre-Trained GloVe Embeddings

Before introducing the attention modules, we describe how the textual input is semantically encoded. Unlike contextualized embeddings such as BERT [44] or GPT [36],

which require task-specific fine-tuning and significant computational resources, static embeddings such as GloVe offer strong semantic priors at minimal computational cost. This makes GloVe an effective choice for attention-based image captioning frameworks [45,46].

(a) Preparation of the Embedding Matrix:

The embedding layer is initialized using pre-trained GloVe vectors from the `glove.6B.200d.txt` file through the following steps:

- Load the vectors into an `embeddings_index` dictionary for efficient retrieval.
- Create an `embedding_matrix` of shape `(vocab_size, embedding_dim)`, initialized with zeros.
- For each vocabulary word, insert its GloVe vector into the matrix; retain the zero vector if the word is not found.

(b) Integration and Semantic Consistency:

The resulting `embedding_matrix` is used to initialize the decoder's embedding layer. To preserve the semantic information encoded in GloVe, the embedding weights are kept non-trainable (frozen) during training. This ensures that the decoder focuses on learning meaningful alignments between visual features and linguistic structure rather than altering established word semantics.

(c) Functionality During Training and Inference:

The embedding layer maps each input token to a dense vector representation, with the source of the input differing depending on the phase:

- **Training phase:** Inputs are ground-truth word indices from the target captions.
- **Inference phase:** Inputs are word indices generated at the previous decoding step.

(d) Role in Attention-Based Decoding:

At each decoding step, the embedding of the current input word is concatenated with the context vector refined by AoA. This joint representation is passed to the GRU unit, which maintains temporal dependencies across the sequence. The linguistic priors encoded in the GloVe embeddings help the model generate fluent, semantically appropriate, and grammatically correct descriptions.

Enhanced Decoder with Bahdanau and Attention on Attention Mechanisms

This subsection details the integration of both Bahdanau Attention and Attention on Attention (AoA) mechanisms within the decoder. These strategies refine spatial and semantic representations, enabling the model to generate more accurate and coherent captions. We begin by describing each mechanism's structure and functionality, followed by an experimental comparison highlighting the added value of AoA over traditional attention.

(a) Bahdanau Attention Mechanism

Bahdanau (additive) attention enables the decoder to adaptively focus on different image regions at each decoding step, thereby improving the contextual alignment between visual features and the generated description.

Mathematical formulation:

$$e_t = \tanh(\mathbf{W}_1 \cdot \mathbf{F} + \mathbf{W}_2 \cdot \mathbf{h}_{t-1}) \quad (1)$$

$$\alpha_t = \text{softmax}(\mathbf{v}^T \cdot e_t) \quad (2)$$

$$\mathbf{c}_t = \sum_i \alpha_{t,i} \mathbf{F}_i \quad (3)$$

Here, \mathbf{F} represents the image features and \mathbf{h}_{t-1} the previous hidden state. The resulting context vector \mathbf{c}_t is concatenated with the embedding of the current input word $\mathbf{E}(w_t)$ and passed to the GRU decoder:

$$\text{GRU_input} = \text{concat}(\mathbf{c}_t, \mathbf{E}(w_t))$$

During training, $\mathbf{E}(w_t)$ corresponds to the ground-truth word; during inference, it refers to the previously generated word. This attention mechanism allows the model to dynamically focus on salient image regions, reinforcing semantic relevance during caption generation.

Nevertheless, Bahdanau attention presents certain limitations. It computes a static attention distribution at each step without internal refinement and considers each region independently, limiting its ability to model relationships among objects. Moreover, condensing the visual information into a single context vector can result in the loss of fine-grained spatial cues, especially in complex scenes.

(b) Attention on Attention (AoA)

To overcome these constraints, we integrate Attention on Attention (AoA), which introduces an additional refinement stage applied to the initial context vector. This layered attention mechanism enables more selective and adaptive focus during decoding, enhancing semantic precision.

Step 1: Primary Attention (Bahdanau)

$$e_t = \tanh(\mathbf{W}_1 \cdot \mathbf{F} + \mathbf{W}_2 \cdot \mathbf{h}_{t-1}) \quad (4)$$

$$\alpha_t = \text{softmax}(\mathbf{v}^T \cdot e_t) \quad (5)$$

$$\mathbf{c}_t = \sum_i \alpha_{t,i} \mathbf{F}_i \quad (6)$$

Step 2: Context Refinement (AoA)

$$\text{score}_2 = \tanh(\mathbf{W}_3 \cdot \mathbf{c}_t) \quad (7)$$

$$\alpha_2 = \text{softmax}(\mathbf{V}_2 \cdot \text{score}_2) \quad (8)$$

$$\mathbf{c}_t^{\text{refined}} = \sum_i \alpha_{2,i} \mathbf{c}_t \quad (9)$$

This refinement brings several benefits:

- **Improved feature selection:** Irrelevant visual cues are more effectively suppressed.
- **Enhanced fluency and coherence:** Captions are more semantically consistent.
- **Better modeling of object interactions:** Particularly useful for complex, multi-object scenes.

The effectiveness of AoA is confirmed by experimental results on the MS COCO dataset (Section 3.2.1, Table 1). Compared to our baseline model, which relies solely on Bahdanau attention, the AoA-enhanced architecture yields notable improvements:

- **BLEU-4:** 38.00 (vs. 0.19)
- **CIDEr:** 128.3 (vs. 12.89)
- **ROUGE:** 66.00 (vs. 21.61)

These results validate the contribution of the dual attention framework, in which the secondary refinement step plays a key role in producing captions that are not only more fluent but also better aligned with the semantic content of the image.

Recurrent Layer (GRU)—Details

The decoder in our image captioning model incorporates a Gated Recurrent Unit (GRU) to model temporal dependencies in conjunction with attention mechanisms. GRUs are known for preserving long-range dependencies while maintaining lower computational

cost compared to LSTMs. Recent findings, such as those by Muškardin et al. [47], demonstrate that GRU hidden states effectively capture semantic progression over time.

1. Input to the GRU Layer

At each decoding step, the GRU processes a combined input that fuses visual and linguistic context. Specifically:

- **Refined context vector:** produced by the AoA mechanism, summarizing salient visual features. Shape: $(\text{batch_size}, \text{embedding_dim})$.
- **Embedded input word:** representing the current token—either a ground-truth word (during training) or a previously generated one (during inference). Shape: $(\text{batch_size}, 1, \text{embedding_dim})$.

To allow concatenation, the context vector is expanded along the temporal axis, resulting in a shape of $(\text{batch_size}, 1, \text{embedding_dim})$. The concatenated input to the GRU is then:

$$(\text{batch_size}, 1, 2 \times \text{embedding_dim})$$

2. GRU Outputs

At each time step t , the GRU produces:

1. **Output vector (Y_t):** a fused representation used by the dense layers for next-word prediction. Shape: $(\text{batch_size}, 1, \text{units})$.
2. **Hidden state (h_t):** a vector of shape $(\text{batch_size}, \text{units})$ that maintains temporal memory throughout caption generation.

Transition from GRU to Dense Layers in Decoder

In our decoder, GRU outputs are transformed into word predictions through two sequential dense layers, D_1 and D_2 , which map the hidden state into a probability distribution over the vocabulary.

Dense Layer D_1 : This layer refines the GRU hidden state \mathbf{H}_t via a non-linear transformation:

$$\mathbf{Z}_t = \phi(\mathbf{H}_t \mathbf{W}_1 + \mathbf{b}_1)$$

where:

- $\mathbf{H}_t \in \mathbb{R}^{N \times d_h}$: GRU hidden state at time step t ;
- $\mathbf{W}_1 \in \mathbb{R}^{d_h \times d_z}$, $\mathbf{b}_1 \in \mathbb{R}^{d_z}$: learnable parameters;
- $\phi(\cdot)$: non-linear activation function (e.g., ReLU);
- $\mathbf{Z}_t \in \mathbb{R}^{N \times d_z}$: refined representation.

Dense Layer D_2 : This layer projects the refined representation into the vocabulary space:

$$\mathbf{O}_t = \mathbf{Z}_t \mathbf{W}_2 + \mathbf{b}_2$$

where:

- $\mathbf{W}_2 \in \mathbb{R}^{d_z \times V}$, $\mathbf{b}_2 \in \mathbb{R}^V$: trainable parameters;
- $\mathbf{O}_t \in \mathbb{R}^{N \times V}$: unnormalized logits over the vocabulary.

These logits are normalized using the softmax function:

$$\hat{\mathbf{Y}}_t = \text{softmax}(\mathbf{O}_t)$$

The predicted word is chosen by:

$$\hat{Y}_t = \arg \max_j \hat{Y}_{t,j}$$

Together, these layers convert the temporal and contextual information encoded in the GRU hidden state into a meaningful probability distribution over the target vocabulary. This transformation enables the decoder to predict the next word in a fluent, grammatically coherent, and contextually grounded manner.

Data Processing Workflow in the Decoder (Summary)

In our implementation, the decoder is defined as a subclass of `tf.keras.Model`, where the `call()` method encapsulates the forward pass logic through the following stages:

1. **Inputs:** The decoder receives:
 - the current token x (ground truth during training, or the previous prediction during inference),
 - the encoded image features from the encoder,
 - and the current hidden state (initialized to zeros at the first decoding step).
2. **Attention Processing:** Bahdanau attention computes an initial context vector by aligning the image features with the decoder's hidden state. This context is then further refined via the Attention on Attention (AoA) module, enhancing focus on the most relevant visual regions.
3. **Embedding and Input Fusion:** The input token x is embedded using a pre-trained GloVe-based embedding layer. The embedding is concatenated with the refined visual context (after expanding its temporal dimension), forming the input to the GRU layer.
4. **GRU and Vocabulary Projection:** The GRU processes this fused input, updating its internal state and generating an output vector. This output is then passed through two dense layers, D_1 and D_2 , which project it into the vocabulary space. Finally, softmax is applied to obtain a probability distribution over possible next words.

The `call()` method thus integrates attention, embedding, sequential modeling, and dense transformations into a unified decoding pipeline, enabling the generation of fluent, coherent, and visually grounded image captions. This method operationalizes all previously described decoder components, combining them into a unified forward pass applicable to both training and inference.

2.4. Implementation Details

To ensure a balance between performance and computational efficiency, our model architecture and hyperparameters were carefully selected.

We use SIFT to extract up to 512 keypoints per image, each described by a 128-dimensional vector. These are clustered via K-Means into 100 visual words ($k = 100$, $n_{\text{init}} = 10$), yielding a compact and expressive local feature representation.

Global features are extracted using a pre-trained ResNet50, which includes 23,587,712 parameters (23,534,592 trainable). For textual input, we utilize 200-dimensional pre-trained GloVe embeddings, offering semantically rich representations with minimal computational cost.

The decoder incorporates a single-layer GRU with 512 units—an architecture widely adopted for its trade-off between learning capacity and efficiency. The Bahdanau attention module projects both 200-dimensional encoder outputs and 512-dimensional hidden states into a shared space for attention computation. This is further refined through an AoA module to enhance focus on salient visual regions.

Training is performed over two stages of 20 epochs each, with a batch size of 64. Each image tensor has shape (64, 2148), merging ResNet features (2048-d) and tiled BoF vectors (100-d). Captions are padded to a uniform length of 51 tokens.

We use the Adam optimizer with a learning rate of 0.001. Data input is streamlined via shuffling (buffer size = 1000) and prefetching (`tf.data.experimental.AUTOTUNE`) to maximize training efficiency.

2.5. Training Strategy

The training strategy is designed to ensure stable convergence, efficient learning, and strong generalization, while optimizing the use of computational resources. Training is conducted in two distinct phases across 40 epochs, and incorporates adaptive checkpointing as well as customized loss management techniques. An Adam optimizer with an initial learning rate of 0.001 is employed to perform efficient gradient-based optimization. A custom Sparse Categorical Cross-Entropy (SCCE) loss function is used to compare integer word indices with predicted probabilities, without incurring the overhead of one-hot encoding.

2.5.1. Handling Padding and Custom Loss Function

To enhance training efficiency, we integrate a masking mechanism that excludes padding tokens from loss computation, preventing the model from learning from irrelevant padding tokens. This ensures that only meaningful tokens contribute to gradient updates, thereby improving model precision. The loss function is defined as:

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T m_t \cdot \log(\hat{y}_{t,y_t})$$

where m_t is a binary mask filtering out padding tokens, ensuring that gradients are only computed for relevant sequences. By assigning a weight of zero to padding tokens (value 0), the model avoids learning unnecessary patterns, which contributes to training stability.

2.5.2. Training and Evaluation Procedures

The training process is managed through two core functions:

- **Training:** The `train_step` function is responsible for optimizing the model's parameters. For each batch, it resets the decoder's hidden state, feeds the ground truth captions word by word (teacher forcing), and computes the loss between the predicted and actual words at each time step. The total loss is then backpropagated to update the model weights. This step does not generate full image descriptions but helps the model learn to do so by improving its internal representations.
- **Evaluation:** The `test_step` function mirrors the training step but without any weight updates. It calculates the loss on the validation set using the model's current parameters, providing a measure of generalization. The validation loss is monitored across epochs, and the model checkpoint with the lowest validation loss is saved automatically to preserve the best-performing version of the model.

This training setup is applied across the full 40-epoch schedule using a two-phase approach with checkpoint-based optimization, as described in the following section.

2.5.3. Two-Phase Training and Checkpoint-Based Performance Monitoring

To ensure both training stability and optimal resource usage, the model is trained in two sequential phases over 40 epochs, with an adaptive checkpointing mechanism that preserves only the best-performing model state. This strategy not only reduces memory overhead, but also improves generalization by focusing on empirically validated checkpoints.

Checkpoints are stored in `Models/CP_caption_ResNet_A0A_256_40k/train`, with the parameter `max_to_keep=1`, ensuring that only the best model is retained throughout the

training process. At the end of each epoch, the validation loss is computed; if it improves, the model weights are saved, guaranteeing that only the most effective model is preserved.

The training is divided into two phases:

- **Phase 1 (Epochs 1–20):** The model is trained from scratch. Checkpoints are saved whenever a new lowest validation loss is observed, ensuring progressive refinement based on actual performance.
- **Phase 2 (Epochs 21–40):** Training resumes from the best checkpoint obtained in Phase 1. This fine-tuning stage allows the model to further optimize its representations while minimizing overfitting risks.

Figure 6 illustrates the evolution of training and validation losses during Phase 1. The consistent decrease in both losses indicates stable convergence and absence of early overfitting.

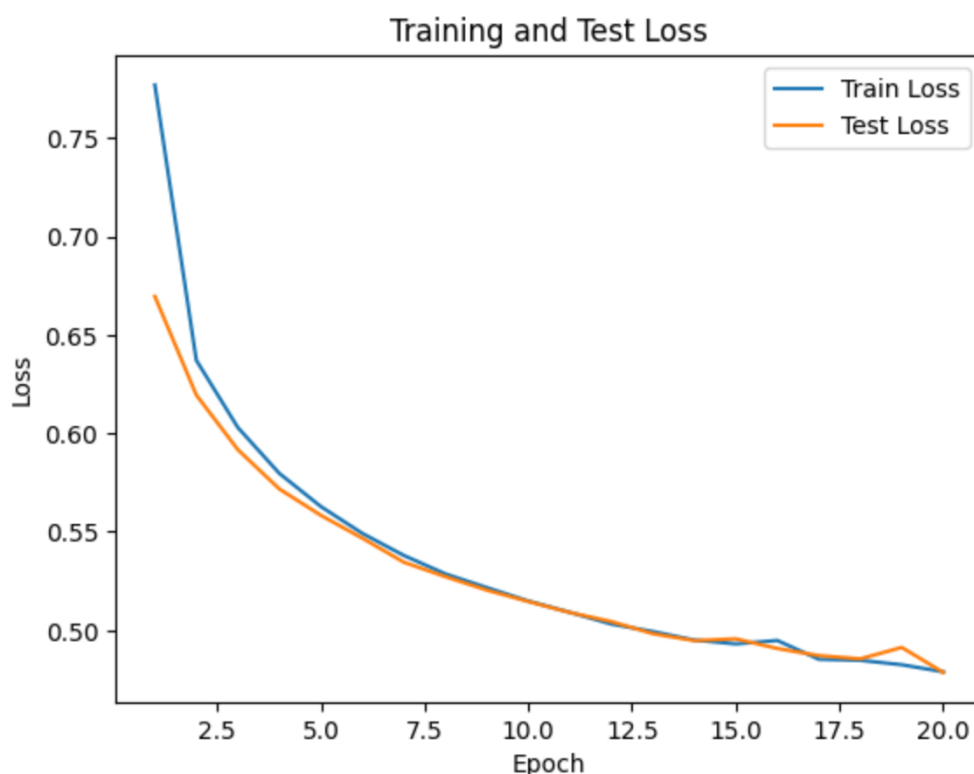


Figure 6. Phase 1: Training and validation losses over the initial 20 epochs. The model exhibits a consistent decline in loss, confirming stable convergence without early signs of overfitting.

In Phase 2 (Figure 7), training continues from the best checkpoint of Phase 1. Despite minor fluctuations between epochs 30 and 35—attributable to input variability—the loss trend remains downward, confirming that the model continues to generalize effectively.

This two-phase training strategy, combined with dynamic checkpointing and careful validation monitoring, ensures efficient training progression while maximizing model generalization performance.

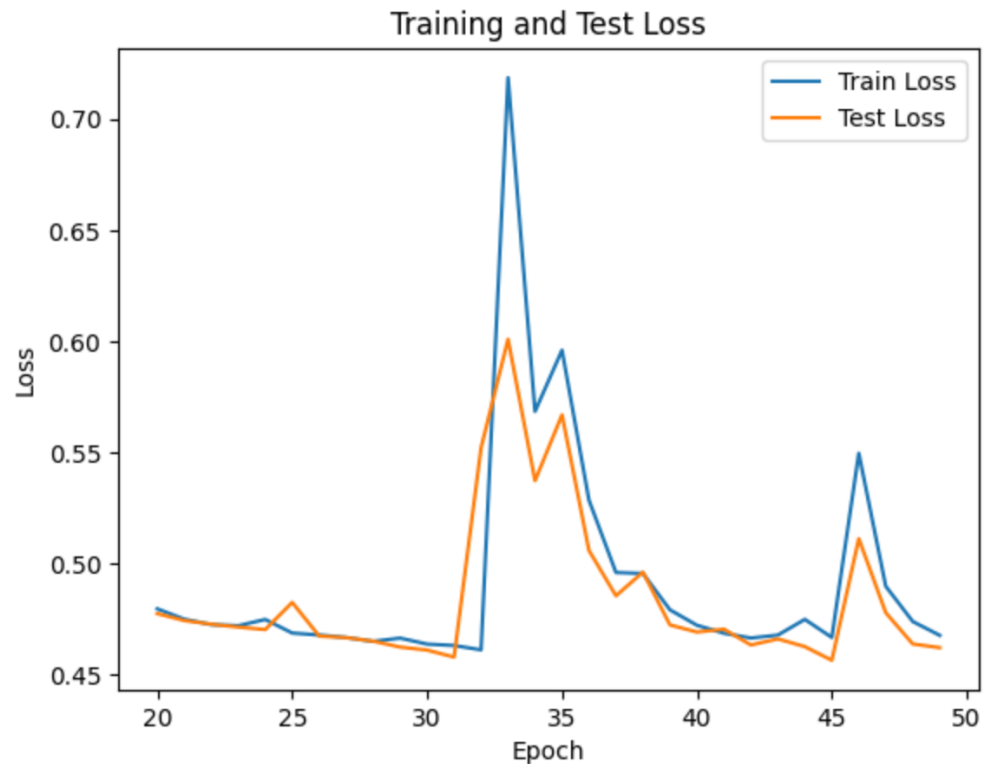


Figure 7. Phase 2: Continued optimization from epochs 21 to 40. The model continues to refine its learning, with minor fluctuations reflecting adaptation to complex input patterns.

2.5.4. Robustness and Error Management

To enhance training robustness, multiple safeguards are integrated:

- **Resource Management:** The phased training strategy distributes computation over two stages, reducing memory overheads and mitigating resource exhaustion.
- **Input Validation:** Data integrity checks are applied to filter out erroneous captions, ensuring high-quality inputs for training.
- **Evaluation Metrics:** Performance metrics such as BLEU [48], METEOR [49], ROUGE [50], CIDEr [51], and SPICE [52] are computed only after validating reference-prediction consistency.
- **Padding Masking:** The custom loss function ensures that training focuses solely on meaningful sequences, improving model precision.

These mechanisms collectively ensure stable training, reduce unnecessary computations, and enhance the model's ability to generalize effectively. The phased approach to training, combined with adaptive checkpointing and a well-structured loss computation method, leads to an optimized learning process while preventing overfitting. Through these strategies, the model maintains high efficiency, leveraging computational resources effectively while achieving strong generalization performance.

2.6. Caption Generation Process

The caption generation process transforms an input image into a descriptive sentence through a structured, step-by-step pipeline that combines classical and deep learning techniques. The `caption_generation(image_path)` function orchestrates this process, integrating both local and global feature extraction, attention mechanisms, and sequential word prediction, as outlined below:

1. **Decoder Initialization:** The decoder's internal state is reset using `reset_state(batch_size=1)` to prepare for caption generation.

2. **Local Feature Extraction (SIFT + K-Means):**
 - `extract_sift_features(image_path)` extracts keypoint descriptors using the SIFT algorithm.
 - K-Means clustering transforms these descriptors into visual words, forming a Bag of Features (BoF) histogram.
3. **Global Feature Extraction (CNN):** High-level semantic features are extracted using a pre-trained CNN via `extract_base_model_features(image_path, new_base_model)`.
4. **Feature Preparation:**
 - Both feature types (BoF and CNN) are reshaped and aligned using `np.tile()` and `np.reshape()`.
 - Features are concatenated using `np.concatenate()` to create a unified representation.
 - This combined feature tensor is fed into the `encoder` for dimensionality reduction.
5. **Caption Decoding:**
 - The decoding process begins with the special start token (`startseq`).
 - The decoder, equipped with Bahdanau Attention, AoA, and a GRU, generates one word at a time.
 - Generation continues until the `endseq` token is predicted.
6. **Output:** The resulting caption is returned either as a list of words or as a TensorFlow tensor.

3. Results

3.1. Dataset Overview and Experimental Setup

To assess the effectiveness of the proposed image captioning framework, experiments were carried out on the MS COCO and Flickr8k datasets. An 80–20 split of MS COCO (train2014) was used for training and testing, while Flickr8k served as an additional benchmark to evaluate generalization capabilities.

All images were resized to **256 × 256 pixels** to ensure uniform input dimensions. For global feature extraction, standard normalization was applied to ResNet-encoded representations. Local descriptors were computed from grayscale versions of the images using SIFT, followed by Bag-of-Features quantization.

The caption generation pipeline was trained in a supervised setting. During evaluation, the model generated captions based solely on visual inputs, without access to ground-truth descriptions. Performance was measured using the metrics detailed in Section 2.5.4.

3.2. Quantitative Performance Comparison Across Datasets

3.2.1. MS COCO Results Analysis

As shown in Table 1, our proposed model achieves top performance on the MS COCO dataset across all reported metrics. It reaches BLEU-1 at 92.04% and CIDEr at 128.3, indicating high fluency and strong semantic alignment in the generated captions.

In particular, our model significantly outperforms the handcrafted-feature baseline, which combines SIFT and Bag of Features for local descriptors, VGG16 for global image encoding, GloVe embeddings for word representation, Bahdanau attention, and a GRU-based language decoder. Compared to this configuration, our model yields substantial improvements: BLEU-1 increases from 12.94% to 92.04%, BLEU-4 from 0.19% to 38.00%, ROUGE from 21.61% to 66.00%, CIDEr from 12.89 to 128.3, and METEOR from 7.34% to 35.00%. These results reflect the impact of incorporating a ResNet backbone, deeper semantic embeddings, and a more refined attention mechanism.

In addition, our model surpasses region-based attention approaches such as Bottom-Up and Top-Down [53], particularly in CIDEr and METEOR, achieving a more effective

balance between lexical precision and contextual richness. This performance gain is attributed to the complementary integration of handcrafted visual descriptors and global deep features, further enhanced by semantically grounded GloVe embeddings.

Table 1. Performance comparison between different image caption generation models on MS COCO.

Dataset	Model	BLEU-1 (%)	BLEU-4 (%)	ROUGE (%)	CIDEr (%)	METEOR (%)
MS COCO	Our Model	92.04	38.00	66.00	128.3	35.00
	SIFT + BoF + VGG16 + GloVe + Bahdanau + GRU (Baseline)	12.94	0.19	21.61	12.89	7.34
	Show, Attend and Tell [25]	71.8	25.0	-	-	23.04
	Attention on Attention [26]	78.7	38.1	58.8	129.8	-
	Bottom-Up and Top-Down [53]	95.2	36.9	57.1	117.9	27.6
	Semantic Attention [54]	91.0	53.4	66.7	168.5	34.1
	Adaptive Attention [55]	74.2	33.2	54.9	108.5	26.6
	GCN-LSTM [56]	80.5	38.2	58.3	127.6	28.5
	ViT-Attn [57]	-	28.69	-	88.79	45.81
	ViT-CNN-Attn [57]	-	31.24	-	95.30	46.98
	Object-Semantic Transformer [58]	80.0	37.7	58.5	132.0	28.9

While Transformer-based models like ViT-CNN-Attn [57] achieve higher METEOR scores—indicating broad semantic coverage—they underperform in BLEU-4 (31.24%) and CIDEr (95.30). This suggests that although their captions reflect global semantic content, they may lack fine-grained sequence-level grounding. Moreover, such models typically demand high-end GPU capacity and large-scale training datasets, which may hinder their adoption in real-world or resource-constrained settings. In contrast, our model maintains competitive performance while operating with a significantly lighter architecture and reduced computational overhead, making it suitable for broader deployment.

3.2.2. Flickr8K Results Analysis

As shown in Table 2, our proposed model achieves the best overall performance on the Flickr8K dataset across all evaluation metrics. It attains the highest n-gram scores, with BLEU-1 at 73.15%, BLEU-2 at 63.02%, BLEU-3 at 47.57%, and BLEU-4 at 37.51%, demonstrating strong fluency and consistent phrase-level accuracy. In terms of semantic evaluation, the model achieves METEOR at 42.98%, CIDEr at 53.00, and SPICE at 20.89%, highlighting effective visual–textual alignment.

Compared to Baseline 1 (SIFT + BoF + VGG16 + Learned Embeddings + Bahdanau + LSTM), our model shows marked gains across all metrics. BLEU-1 improves from 59.57% to 73.15%, BLEU-2 from 47.24% to 63.02%, BLEU-3 from 39.33% to 47.57%, and BLEU-4 from 33.78% to 37.51%. Semantic scores also rise significantly, with METEOR increasing from 20.34% to 42.98%, CIDEr from 17.98 to 53.00, and SPICE from 11.63% to 20.89%. These improvements illustrate the impact of our attention-enhanced dual-fusion strategy and the integration of deeper contextual visual features.

Relative to Baseline 2 (SIFT + BoF + VGG16 + GloVe Embeddings + Bahdanau + LSTM), our model further demonstrates its advantage, particularly in semantic metrics. METEOR improves from 20.94% to 42.98%, CIDEr from 0.17 to 53.00, and SPICE from 12.75% to

20.89%. These results confirm the value of deeper contextual representation and robust visual–semantic fusion.

Our model also surpasses the recent method by Al Badarneh et al. (2023) in both METEOR and SPICE, while remaining competitive in CIDEr. Additionally, earlier models such as Google NIC, Log-Bilinear, and Soft-Attention consistently report lower scores across all metrics, reaffirming the effectiveness of our multi-scale fusion and refined attention mechanisms.

Table 2. Performance comparison between different image caption generation models on Flickr8K.

Dataset	Model	BLEU-1 (%)	BLEU-2 (%)	BLEU-3 (%)	BLEU-4 (%)	ROUGE-L (%)	METEOR (%)	CIDEr (%)	SPICE (%)
Flickr8K	Our Model (Current Approach)	73.15	63.02	47.57	37.51	43.60	42.98	53.00	20.89
	Our Baseline 1: SIFT + VGG16 + BoF + Learned Embeddings + Bahdanau + LSTM	59.57	47.24	39.33	33.78	43.54	20.34	17.98	11.63
	Our Baseline 2: SIFT + BoF + VGG16 + GloVe Embeddings + Bahdanau + LSTM	59.42	0.61	0.11	0.04	43.63	20.94	0.17	12.75
	Al Badarneh et al. (2023) [59]	72.80	49.50	32.30	20.80	43.20	23.50	60.40	16.40
	Jiang et al. (2019) [60]	69.00	47.10	32.40	21.90	50.20	20.30	50.70	-
	Patel et al. (2020) [61]	60.10	41.40	27.40	18.10	43.30	18.30	45.20	-
	Google NIC [19]	63.00	41.00	27.00	-	-	-	-	-
	Log Bilinear [12]	65.60	42.40	27.70	17.70	-	17.31	-	-
	Soft-Attention [25]	67.00	44.80	29.90	19.50	-	18.93	-	-
	Hard-Attention [25]	67.00	45.70	31.40	21.30	-	20.30	-	-
CF [38]	64.50	46.70	32.70	22.80	-	21.30	-	-	

3.3. Qualitative Evaluation of Caption Generation on the MS COCO and Flickr8K Datasets

This section provides a qualitative analysis of the captions generated by our model on selected examples from the MS COCO and Flickr8K datasets. For additional outputs, Appendix A presents a broader range of captioning results produced by our model on both datasets.

Moreover, Appendix B provides a side-by-side qualitative comparison between captions generated by our proposed model and a baseline approach based on handcrafted and deep features (SIFT + BoF + VGG16), GloVe word embeddings, Bahdanau attention, and a GRU decoder. This comparison further illustrates the improvements in fluency, semantic precision, and visual grounding achieved by our current model.

Table 3 illustrates several captions generated by our model on the MS COCO dataset. These examples highlight the model’s ability to produce accurate and contextually relevant descriptions by effectively grounding the generated text in the visual content.

For instance, the caption “A large stone clock tower” provides a concise and correct semantic summary, focusing on the salient object with clarity. “A Delta airline taking off from the airport” demonstrates the model’s capacity to recognize branded objects and dynamic actions, which are often challenging in vision-language alignment tasks.

The output “A man in a suit and tie standing next to a large group of people” captures both the subject and a relational context, illustrating the model’s understanding of human presence and grouping. Similarly, “A man and a woman standing in a living room” reflects a grasp of indoor environments and social scenarios.

Although the caption “An elephant standing next to a rock, directing something” contains a slight ambiguity due to the verb “directing,” it still successfully identifies the central object and its spatial relationship. Such occasional semantic vagueness is rare and

does not undermine the overall quality of the generated captions. Instead, it highlights a potential avenue for future enhancement in fine-grained activity recognition.

Overall, these examples confirm the model's capacity to generate syntactically correct, semantically rich, and visually grounded descriptions in diverse and complex scenes present in the MS COCO dataset.

Table 3. Caption Generation Results on the MS COCO Dataset.

Image	Generated Caption
	A large stone clock tower.
	A Delta airline taking off from the airport.
	A man in a suit and tie standing next to a large group of people.
	A man and a woman standing in a living room.
	An elephant standing next to a rock, directing something.

We now present additional qualitative results on the Flickr8K dataset to further assess the model's performance on a wider variety of informal and diverse visual scenes.

Table 4 shows several image-caption examples drawn from the Flickr8K dataset. The results confirm the model's ability to generate context-aware, semantically rich, and fluent descriptions across a wide range of everyday scenarios.

Captions such as “The dog is standing in front of a fence in a grassy yard” and “A young child sitting in a car” reflect precise object recognition and spatial understanding. Similarly, the caption “A sheep is standing on a lush green field” demonstrates a coherent understanding of the relationship between subject and background, reinforcing the model’s scene comprehension.

However, some minor linguistic imperfections can still be observed. For example, in the caption “A group of people are standing around a long wooden table with a group of people watch” the verb “watch” is grammatically incorrect and the sentence structure is somewhat redundant. Despite this, the scene is still well understood in terms of its core elements (people, table, social interaction), indicating that the model has correctly grounded the main visual components. This example illustrates how the system can capture the essence of a scene even when the output lacks full grammatical precision.

Although such grammatical inconsistencies are rare, they offer useful insights for refining linguistic post-processing and improving action disambiguation. Overall, these qualitative examples validate the robustness and flexibility of our model across both structured (MS COCO) and informal (Flickr8K) datasets. The integration of handcrafted features (SIFT + BoF) with deep visual embeddings and an enhanced attention mechanism plays a key role in generating coherent and contextually grounded descriptions. Future work may focus on reducing redundancy and further enhancing language fluency. A side-by-side qualitative comparison between the proposed model and a baseline architecture is provided in Appendix B.

Table 4. Caption Generation Results on the Flickr8K Dataset.






Image	Generated Caption
	The dog is standing in front of a fence in a grassy yard.
	A group of people in snow skis.
	A group of people are standing around a long wooden table with a group of people watch.
	A young child sitting in a car.

Table 4. Cont.

Image	Generated Caption
	A sheep is standing on a lush green field.

4. Discussion and Conclusions

In this study, we presented a hybrid image captioning model that integrates both hand-crafted and deep learning-based visual features to enhance semantic understanding and linguistic fluency in generated captions. Specifically, we combined local descriptors—SIFT and Bag of Features (BoF)—with high-level global representations extracted via ResNet50. This complementary fusion enables the system to capture both fine-grained textures and holistic scene semantics.

The decoder incorporates Bahdanau Attention and the Attention-on-Attention (AoA) mechanism, allowing dynamic focus on salient visual regions. GRU-based sequence modeling, enhanced by GloVe embeddings, contributes to syntactic coherence. The model achieves competitive performance on MS COCO and Flickr8K datasets, with a CIDEr score of 128.3 and strong BLEU and METEOR scores, validating its effectiveness over conventional baselines.

Despite its strengths, several aspects of the hybrid Bahdanau+AoA architecture present avenues for enhancement:

- **Long-Range Dependency Modeling:** GRUs are effective for local context but less suitable for capturing global dependencies. Replacing the decoder with a Transformer-based module may improve handling of complex spatial and semantic relationships.
- **Attention Error Propagation:** Misaligned initial attention from Bahdanau may propagate through AoA. Incorporating cross-modal regularization or uncertainty-aware mechanisms could help mitigate this.
- **Feature Fusion Redundancy:** Concatenating high-dimensional ResNet and BoF features can introduce redundancy. Adaptive gating or dimensionality reduction strategies may optimize this fusion.
- **Generic or Repetitive Captions:** Some outputs are overly generic due to dataset biases or decoder limitations. Techniques like diverse beam search or adversarial training could foster output diversity.
- **Inference Overhead:** The dual-attention mechanism introduces computational cost. Lightweight alternatives, such as attention pruning or knowledge distillation, could improve deployment feasibility.

These limitations offer valuable insights for future innovation. Addressing them will enable more robust, scalable, and semantically precise captioning systems.

Future research could explore Transformer-based visual encoders (e.g., ViT, Swin Transformer [62]) and contextual language models (e.g., BERT, GPT) to further improve visual-textual alignment. Reinforcement learning approaches such as self-critical sequence training (SCST), combined with large-scale multimodal datasets, may also enhance generalization and real-world applicability.

By bridging classical descriptors with modern attention-driven decoding, our hybrid model provides a robust and interpretable framework for image captioning. The results across benchmarks confirm its potential, particularly in resource-constrained or diverse-

data environments. This work lays the foundation for more expressive and human-like caption generation, with promising applications in assistive technologies, automated content description, and beyond.

Author Contributions: Conceptualization, C.S.; methodology, C.S. and M.R.; software, C.S.; validation, C.S. and I.B.; writing—original draft preparation, C.S.; writing—review and editing, C.S. and I.B.; supervision, M.R. All authors have read and agreed to the published version of the manuscript.

Funding: This research did not receive any funding from public, commercial, or non-profit organizations.

Data Availability Statement: The datasets generated and analyzed during this study are available from the corresponding author upon reasonable request.

Conflicts of Interest: The authors declare no conflicts of interest.

Appendix A. Additional Captions Generated by the Proposed Method: Unseen Examples from the MS COCO and Flickr8K Datasets

Appendix A.1. Generated Captions for Unseen Images from the MS COCO Test Set



Figure A1. Generated caption: "A man riding a bike in the city".



Figure A2. Generated caption: "A display of freshly made donuts, muffins, and snacks".



Figure A3. Generated caption: "A stop sign is on a street corner".



Figure A4. Generated caption: "A man and a woman eat sandwiches while a dog watches".

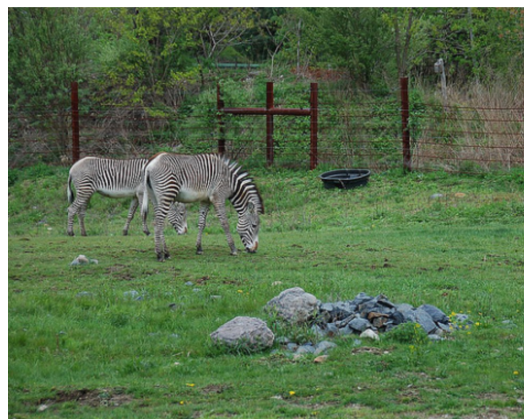


Figure A5. Generated caption: "Zebras standing next to each other".



Figure A6. Generated caption: "A man is standing on a skateboard in front of a parking area".



Figure A7. Generated caption: "A group of women talking around a wine bar".



Figure A8. Generated caption: "A cow standing in a field under the sky".

Appendix A.2. Generated Captions for Unseen Images from the Flickr8K Test Set



Figure A9. Generated caption: "A white dog carrying a frisbee in a muddy area".



Figure A10. Generated caption: "A couple sitting by a lake with their baby".



Figure A11. Generated caption: "A dog is standing in a grassy area".



Figure A12. Generated caption: "A brown dog is walking along the beach".



Figure A13. Generated caption: "A young boy is jumping on a couch".



Figure A14. Generated caption: "A brown dog pursues a frisbee across the grass".



Figure A15. Generated caption: "A group of boats are docked in a marina".

Appendix B. Comparison of Generated Captions (Our Model vs. Baseline: SIFT + BoF + VGG16 + GloVe + Bahdanau + GRU)

The following examples illustrate differences between captions generated by our proposed architecture and a baseline model. The baseline combines handcrafted and deep visual features using SIFT, Bag of Features (BoF), and VGG16 for image encoding. It employs a decoder that integrates GloVe word embeddings, Bahdanau attention, and a GRU-based language model.

This configuration corresponds to the baseline evaluated in Section 3.2.1 (Table 1). To improve upon this architecture, our approach incorporates a dual-branch visual encoding strategy alongside a refined attention mechanism based on Attention on Attention (AoA). The following qualitative examples illustrate how these enhancements lead to more accurate, fluent, and semantically grounded captions.

Table A1. Comparison of Generated Captions on Flickr8K (Baseline vs. Our Approach) *Baseline:* SIFT + BoF + VGG16 + GloVe + Bahdanau + GRU.





Image	Baseline	Our Approach
	puppy figure.	The dog is standing in front of a fence in a grassy yard.
	a man stands in snow covered mountain house.	A group of people in snow skis.
	a young child is sitting outdoors.	A young child sitting in a car.
	a group of kids look on a couch.	A group of people are standing around a long wooden table with a group of people watch.

Table A2. Comparison of Generated Captions on the MS COCO Dataset (Baseline vs. Our Approach) *Baseline:* SIFT + BoF + VGG16 + GloVe + Bahdanau + GRU.



Image	Baseline	Our Approach
	a large tower with a clock tower	A large stone clock tower

Table A2. Cont.

Image	Baseline	Our Approach
	a man and woman staining looking at their phones	A couple of people standing in front of a wine glass

References

- Gurari, D.; Zhao, Y.; Zhang, M.; Bhattacharya, N.; Bigham, J.P.; Grauman, K. Captioning Images Taken by People Who Are Blind. In Proceedings of the European Conference on Computer Vision (ECCV), Glasgow, UK, 23–28 August 2020; pp. 417–434. [CrossRef]
- Geiger, A.; Lenz, P.; Urtasun, R. Are We Ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Providence, RI, USA, 16–21 June 2012; pp. 3354–3361. [CrossRef]
- Hossain, M.Z.; Sohel, F.; Shiratuddin, M.F.; Laga, H. A Comprehensive Survey of Deep Learning for Image Captioning. *ACM Comput. Surv. (CSUR)* **2019**, *51*, 1–36. [CrossRef]
- Mallat, S. A Theory for Multiresolution Signal Decomposition: The Wavelet Representation. *IEEE Trans. Pattern Anal. Mach. Intell.* **1989**, *11*, 674–693. [CrossRef]
- Daubechies, I. *Ten Lectures on Wavelets*; Society for Industrial and Applied Mathematics (SIAM): Philadelphia, PA, USA, 1992. [CrossRef]
- Rabiner, L.R. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proc. IEEE* **1989**, *77*, 257–286. [CrossRef]
- Lowe, D.G. Object Recognition from Local Scale-Invariant Features. In Proceedings of the Seventh IEEE International Conference on Computer Vision (ICCV), Corfu, Greece, 20–25 September 1999; pp. 1150–1157. [CrossRef]
- Csurka, G.; Dance, C.R.; Fan, L.; Willamowski, J.; Bray, C. Visual Categorization with Bags of Keypoints. In Proceedings of the ECCV Workshop on Statistical Learning in Computer Vision, Prague, Czech Republic, 16 May 2004; pp. 1–22.
- Farhadi, A.; Hejrati, M.; Sadeghi, M.A.; Young, P.; Rashtchian, C.; Hockenmaier, J.; Forsyth, D. Every Picture Tells a Story: Generating Sentences from Images. In Proceedings of the 11th European Conference on Computer Vision (ECCV 2010), Heraklion, Crete, Greece, 5–11 September 2010; Springer: Berlin/Heidelberg, Germany, 2010; pp. 15–29. [CrossRef]
- Kulkarni, G.; Premraj, V.; Ordonez, V.; Dhar, S.; Li, S.; Choi, Y.; Berg, A.C.; Berg, T.L. BabyTalk: Understanding and Generating Simple Image Descriptions. *IEEE Trans. Pattern Anal. Mach. Intell.* **2013**, *35*, 2891–2903. [CrossRef]
- Mnih, A.; Hinton, G.E. Three New Graphical Models for Statistical Language Modelling. In Proceedings of the 24th International Conference on Machine Learning (ICML 2007); ACM: New York, NY, USA, 2007; pp. 641–648. [CrossRef]
- Kiros, R.; Salakhutdinov, R.; Zemel, R.S. Multimodal Neural Language Models. In Proceedings of the 31st International Conference on Machine Learning (ICML), Beijing, China, 21–26 June 2014; pp. II-595–II-603. Available online: <https://dl.acm.org/doi/10.5555/3044805.3044959> (accessed on 19 March 2025).
- Kiros, R.; Salakhutdinov, R.; Zemel, R. Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models. *arXiv* **2014**, arXiv:1411.2539.
- LeCun, Y.; Bottou, L.; Bengio, Y.; Haffner, P. Gradient-Based Learning Applied to Document Recognition. *Proc. IEEE* **1998**, *86*, 2278–2324. [CrossRef]
- Hochreiter, S.; Schmidhuber, J. Long Short-Term Memory. *Neural Comput.* **1997**, *9*, 1735–1780. [CrossRef]
- Cho, K.; van Merriënboer, B.; Gulcehre, C.; Bahdanau, D.; Bougares, F.; Schwenk, H.; Bengio, Y. Learning Phrase Representations Using RNN Encoder–Decoder for Statistical Machine Translation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Doha, Qatar, 25–29 October 2014; pp. 1724–1734. Available online: <https://aclanthology.org/D14-1179> (accessed on 19 March 2025).

17. Mikolov, T.; Karafiát, M.; Burget, L.; Černocký, J.; Khudanpur, S. Recurrent Neural Network Based Language Model. *Proc. Interspeech* **2010**, *2*, 1045–1048. Available online: https://www.isca-speech.org/archive/interspeech_2010/mikolov10_interspeech.html (accessed on 19 March 2025).
18. Mao, J.; Xu, W.; Yang, Y.; Wang, J.; Huang, Z.; Yuille, A. Learning Like a Child: Fast Novel Visual Concept Learning from Sentence Descriptions of Images. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), Santiago, Chile, 13–16 December 2015; pp. 2533–2541. Available online: https://openaccess.thecvf.com/content_iccv_2015/papers/Mao_Learning_Like_a_ICCV_2015_paper.pdf (accessed on 19 March 2025).
19. Vinyals, O.; Toshev, A.; Bengio, S.; Erhan, D. Show and Tell: A Neural Image Caption Generator. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Boston, MA, USA, 7–12 June 2015; pp. 3156–3164. [CrossRef]
20. Girshick, R.; Donahue, J.; Darrell, T.; Malik, J. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Columbus, OH, USA, 23–28 June 2014; pp. 580–587. [CrossRef]
21. Girshick, R. Fast R-CNN. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), Santiago, Chile, 7–13 December 2015; pp. 1440–1448. [CrossRef]
22. Ren, S.; He, K.; Girshick, R.; Sun, J. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In Proceedings of the Advances in Neural Information Processing Systems 28 (NeurIPS 2015), Montreal, QC, Canada, 7–12 December 2015; pp. 91–99. Available online: https://proceedings.neurips.cc/paper_files/paper/2015/file/14bfa6bb14875e45bba028a21ed38046-Paper.pdf (accessed on 19 March 2025).
23. Redmon, J.; Divvala, S.; Girshick, R.; Farhadi, A. You Only Look Once: Unified, Real-Time Object Detection. In Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, NV, USA, 27–30 June 2016; pp. 779–788. [CrossRef]
24. Karpathy, A.; Fei-Fei, L. Deep Visual-Semantic Alignments for Generating Image Descriptions. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Boston, MA, USA, 7–12 June 2015; pp. 3128–3137.
25. Xu, K.; Ba, J.; Kiros, R.; Cho, K.; Courville, A.; Salakhudinov, R.; Zemel, R.; Bengio, Y. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In Proceedings of the International Conference on Machine Learning (ICML), Lille, France, 6–11 July 2015; pp. 2048–2057.
26. Huang, L.; Wang, W.; Chen, J.; Wei, X.Y. Attention on Attention for Image Captioning. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), Seoul, Republic of Korea, 27 October–2 November 2019; pp. 4634–4643. [CrossRef]
27. Al-Malla, M.A.; Jafar, A.; Ghneim, N. Image Captioning Model Using Attention and Object Features to Mimic Human Image Understanding. *J. Big Data* **2022**, *9*, 20. [CrossRef]
28. Fang, H.; Deng, L.; Mitchell, M.; Gupta, S.; Dollár, P.; Platt, J.C.; Iandola, F.; Gao, J.; Zitnick, C.L.; Srivastava, R.K.; et al. From Captions to Visual Concepts and Back. *arXiv* **2015**, arXiv:1411.4952v3.
29. Yang, X.; Wang, Y.; Wang, Z.; Xu, Y.; Bai, X.; Bai, S. Auto-Encoding Scene Graphs for Image Captioning. In Proceedings of the 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA, 15–20 June 2019; pp. 10685–10694. [CrossRef]
30. Xu, D.; Zhu, Y.; Choy, C.B.; Fei-Fei, L. Scene Graph Generation by Iterative Message Passing. In Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, USA, 21–26 July 2017; pp. 5410–5419. [CrossRef]
31. He, K.; Gkioxari, G.; Dollár, P.; Girshick, R. Mask R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Venice, Italy, 22–29 October 2017; pp. 2961–2969. [CrossRef]
32. Felzenszwalb, P.F.; Huttenlocher, D.P. Efficient Graph-Based Image Segmentation. *Int. J. Comput. Vis.* **2004**, *59*, 167–181. [CrossRef]
33. Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. An Image Is Worth 16 × 16 Words: Transformers for Image Recognition at Scale. *arXiv* **2021**, arXiv:2010.11929. Available online: <https://arxiv.org/abs/2010.11929> (accessed on 19 March 2025).
34. Liu, W.; Chen, S.; Guo, L.; Zhu, X.; Liu, J. CPTN: Full Transformer Network for Image Captioning. *arXiv* **2021**, arXiv:2101.10804. Available online: <https://arxiv.org/abs/2101.10804> (accessed on 19 March 2025).
35. Radford, A.; Kim, J.W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. Learning Transferable Visual Models From Natural Language Supervision. In Proceedings of the International Conference on Machine Learning (ICML), Virtual, 18–24 July 2021. Available online: <https://arxiv.org/abs/2103.00020> (accessed on 19 March 2025).
36. OpenAI. GPT-4V(ision): Extending GPT-4 to Multimodal Inputs. OpenAI Technical Report, 2023. Available online: <https://cdn.openai.com/papers/gpt-4.pdf> (accessed on 12 October 2023).
37. Alayrac, J.B.; Donahue, J.; Luc, P.; Miech, A.; Barr, I.; Hasson, Y.; Lenc, K.; Mensch, A.; Millican, K.; Reynolds, M.; et al. Flamingo: A Visual Language Model for Few-Shot Learning. *Adv. Neural Inf. Process. Syst.* **2022**, *35*, 23716–23736.
38. Kalimuthu, M.; Mogadala, A.; Mosbach, M.; Klakow, D. Fusion Models for Improved Visual Captioning. *arXiv* **2020**, arXiv:2010.15251.

39. Pennington, J.; Socher, R.; Manning, C.D. GloVe: Global Vectors for Word Representation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Doha, Qatar, 25–29 October 2014; Association for Computational Linguistics: Stroudsburg, PA, USA, 2014; pp. 1532–1543. Available online: <https://aclanthology.org/D14-1162> (accessed on 19 March 2025).
40. He, K.; Zhang, X.; Ren, S.; Sun, J. Deep Residual Learning for Image Recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, NV, USA, 27–30 June 2016; pp. 770–778. [CrossRef]
41. MacQueen, J.B. Some Methods for Classification and Analysis of Multivariate Observations. In Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Berkeley, CA, USA, 21 June–18 July 1965 and 27 December 1965–7 January 1966; University of California Press: Oakland, CA, USA, 1967; Volume 1, pp. 281–297.
42. Simonyan, K.; Zisserman, A. Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv* **2014**, arXiv:1409.1556. Available online: <https://arxiv.org/abs/1409.1556> (accessed on 12 July 2021).
43. Maas, A.L.; Hannun, A.Y.; Ng, A.Y. Rectifier Nonlinearities Improve Neural Network Acoustic Models. In Proceedings of the 30th International Conference on Machine Learning (ICML 2013), Workshop on Deep Learning for Audio, Speech and Language Processing, Atlanta, GA, USA, 16–21 June 2013. Available online: https://ai.stanford.edu/~amaas/papers/relu_hybrid_icml2013_final.pdf (accessed on 19 March 2025).
44. Devlin, J.; Chang, M.-W.; Lee, K.; Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT), Minneapolis, MN, USA, 2–7 June 2019; Association for Computational Linguistics: Stroudsburg, PA, USA, 2019; pp. 4171–4186. Available online: <https://aclanthology.org/N19-1423> (accessed on 19 March 2025).
45. Quadrai. Image Captioning with Attention, ResNet and GloVe. Available online: <https://www.kaggle.com/code/quadrai/image-captioning-with-attention-resnet-and-glove> (accessed on 19 March 2025).
46. Anundskås, L.H.; Afridi, H.; Tarekegn, A.N.; Yamin, M.M.; Ullah, M.; Yamin, S.; Cheikh, F.A. GloVe-Ing Attention: A Multi-Modal Neural Learning Approach to Image Captioning. In Proceedings of the 2023 IEEE International Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW), Rhodes Island, Greece, 4–9 June 2023. [CrossRef]
47. Muškardin, E.; Tappler, M.; Pill, I.; Aichernig, B.K.; Pock, T. On the Relationship Between RNN Hidden-State Vectors and Semantic Structures. In *Findings of the Association for Computational Linguistics: ACL 2024*; Association for Computational Linguistics: Bangkok, Thailand, 2024; pp. 5641–5658. Available online: <https://aclanthology.org/2024.findings-acl.335/> (accessed on 22 February 2024).
48. Papineni, K.; Roukos, S.; Ward, T.; Zhu, W.-J. BLEU: A Method for Automatic Evaluation of Machine Translation. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL), Philadelphia, PA, USA, 6–12 July 2002; pp. 311–318. [CrossRef]
49. Banerjee, S.; Lavie, A. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures*; Association for Computational Linguistics: Ann Arbor, MI, USA, 2005; pp. 65–72.
50. Lin, C.Y. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out: Proceedings of the ACL Workshop*; Association for Computational Linguistics: Ann Arbor, MI, USA, 2004; pp. 74–81.
51. Vedantam, R.; Lawrence Zitnick, C.; Parikh, D. CIDEr: Consensus-Based Image Description Evaluation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Boston, MA, USA, 7–12 June 2015; pp. 4566–4575. [CrossRef]
52. Anderson, P.; Fernando, B.; Johnson, M.; Gould, S. SPICE: Semantic Propositional Image Caption Evaluation. In *European Conference on Computer Vision (ECCV)*; Springer: Cham, Switzerland, 2016; pp. 382–398. [CrossRef]
53. Anderson, P.; He, X.; Buehler, C.; Teney, D.; Johnson, M.; Gould, S.; Zhang, L. Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Salt Lake City, UT, USA, 18–22 June 2018; pp. 6077–6086. Available online: https://openaccess.thecvf.com/content_cvpr_2018/papers/Anderson_Bottom-Up_and_Top-Down_CVPR_2018_paper.pdf (accessed on 15 July 2023).
54. You, Q.; Jin, H.; Wang, Z.; Fang, C.; Luo, J. Image Captioning with Semantic Attention. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, NV, USA, 26 June–1 July 2016; pp. 4651–4659.
55. Lu, J.; Xiong, C.; Parikh, D.; Socher, R. Knowing When to Look: Adaptive Attention via a Visual Sentinel for Image Captioning. *arXiv* **2017**, arXiv:1612.01887.
56. Yao, T.; Pan, Y.; Li, Y.; Mei, T. Exploring Visual Relationship for Image Captioning. *arXiv* **2018**, arXiv:1809.07041.
57. Cahyono, J.A.; Jusuf, J.N. Automated Image Captioning with CNNs and Transformers. *arXiv* **2024**, arXiv:2412.10511.
58. Hafeth, D.A.; Kollias, S. Insights into Object Semantics: Leveraging Transformer Networks for Advanced Image Captioning. *Sensors* **2024**, *24*, 1796. [CrossRef]

59. Al Badarneh, I.; Hammo, B.H.; Al-Kadi, O. An Ensemble Model with Attention-Based Mechanism for Image Captioning. *Information* **2023**, *14*, 56. Available online: <https://www.mdpi.com/2078-2489/14/1/56> (accessed on 10 October 2024). [[CrossRef](#)]
60. Jiang, L.; Zhang, Z.; Huang, Z.; Tan, T. Incorporating Copying Mechanism in Image Captioning for Learning Novel Objects. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA, 15–20 June 2019; pp. 5263–5271. Available online: https://openaccess.thecvf.com/content_cvpr_2017/papers/Yao_Incorporating_Copying_Mechanism_CVPR_2017_paper.pdf (accessed on 10 October 2024).
61. Patel, A.; Lakhotia, K.; Jain, N.; Bhattacharyya, C. Diversity-Promoting GAN for Generating Natural Image Descriptions. In Proceedings of the AAAI Conference on Artificial Intelligence, New York, NY, USA, 7–12 February 2020; Volume 34, pp. 11635–11642. Available online: <https://ojs.aaai.org/index.php/AAAI/article/view/6800> (accessed on 10 October 2024).
62. Liu, Z.; Lin, Y.; Cao, Y.; Hu, H.; Wei, Y.; Zhang, Z.; Lin, S.; Guo, B. Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), Montreal, QC, Canada, 11–17 October 2021; pp. 10012–10022. Available online: https://openaccess.thecvf.com/content/ICCV2021/papers/Liu_Swin_Transformer_Hierarchical_Vision_Transformer_Using_Shifted_Windows_ICCV_2021_paper.pdf (accessed on 19 March 2025).

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.