



A Practical CPU-Friendly Model for Automated Image Caption Generation

Prof. Ashwini P¹, Mohammed Shadab Khalid², Mariam Fatima³, Rifa Habeeb⁴, Sonali Naag⁵

^{1,2,3,4,5}Department of Computer Science and Engineering, ATME College of Engineering, Mysore, Karnataka, India.

ABSTRACT

Automated image captioning aims to convert visual information into meaningful textual descriptions. However, many existing approaches rely on computationally heavy architectures that are unsuitable for real-time use on standard hardware. This paper presents a CPU-friendly and lightweight model that generates captions efficiently without requiring high-end resources. The system uses a compact visual feature extractor and a streamlined text-generation module capable of producing grammatically sound and context-aware descriptions. A user-oriented web interface is implemented using Streamlit for real-time interaction. Experimental evaluations performed on a custom dataset demonstrate strong relevance, low latency, and high usability. The proposed approach is practical for accessibility tools, academic environments, and lightweight applications where efficiency and simplicity are essential.

Keywords—Image Captioning, Visual Understanding, Natural Language Generation, Lightweight Model, Real-Time Processing.

INTRODUCTION

With the rapid increase in digital images across the internet, mobile devices, and surveillance systems, automating the description of images has become an important task. Automatic image caption generation sits at the intersection of computer vision and natural language processing, aiming to replicate the human ability to understand visual scenes and express them through natural language. Many modern captioning systems depend on large, resource-intensive architectures that limit their practical use in low-resource or academic environments. The motivation behind this project is to design a caption generator that works smoothly on regular CPUs, requires minimal computational power, produces relevant and understandable captions, and operates through a simple web-based interface. This paper presents such a solution, built to be lightweight, accessible, and real-time.

II. PROBLEM STATEMENT

Current image captioning solutions often require high-end GPUs, heavy memory consumption, complex deployment steps, and long inference time. As a result, they are not suitable for students, small organizations, academic projects, or environments with limited hardware. This study addresses the lack of efficient and CPU-compatible captioning systems that provide meaningful output while maintaining low computational load.

III. OBJECTIVES

- To design a compact visual feature extraction module.
- To generate meaningful and coherent captions.
- To build a real-time, user-friendly interface.
- To evaluate caption quality using small-scale metrics.
- To maintain low processing time and high usability.
- To create a deployable system suitable for CPU-only environments.

IV. LITERATURE REVIEW

Image captioning has come a long way since the days of pattern matching. It has grown to use deep learning-based vision-language architectures. This section will talk about some important early and modern works in this area.

- Meshed-Memory Transformer (CVPR 2020)

Cornia et al. introduced a memory-augmented transformer for captioning. Limitation: GPU-dependent and computationally heavy. Improvement in our work: Light and fast CPU-only system.

- Unified Vision-Language Models (ECCV 2021)

Pretrained multimodal models perform well on captioning tasks. Limitation: Massive parameter sizes. Our improvement: Smaller model and faster inference.

- CLIP-Guided Captioning (2022)

Uses CLIP embeddings for context alignment. Limitation: Slow on CPUs due to heavy encoders. Our improvement: Lightweight visual feature extraction.

- Mobile-Friendly Captioning Models (2023)

Focus on compressed models for mobile hardware. Limitation: Accuracy drops significantly. Our improvement: Balanced accuracy and speed.

- Assistive Technology Captioning (2024)

Developed for visually impaired assistance. Limitation: Prototype-level systems. Our improvement: Fully deployable and user-ready model.

- Efficient Scene Description (2025)

Research into hybrid visual-semantic models is ongoing. Limitation: Reliance on GPU acceleration. Our improvement: A stable CPU-only implementation.

V. PROPOSED SYSTEM

The proposed system is designed around three key principles:

- Lightweight processing
- Real-time caption generation
- User-friendly interaction

VI. APPLICATIONS

- Tools for visually impaired users
- Social media automation
- E-commerce product description
- Smart surveillance
- Education and research projects
- Digital asset management

VII. EVALUATION

Caption quality was evaluated using a custom dataset and standard metrics:

- BLEU Score
- METEOR Score
- ROUGE-L Score

Results demonstrate strong semantic alignment and low caption latency despite CPU-only inference.

VIII. METHODOLOGY

A. Image Preprocessing

- Image resizing

- Normalization
- Conversion to model-friendly format

These steps ensure uniform processing and faster computation.

B. Visual Feature Extraction

A compact vision module extracts:

- Edges
- Objects
- Shapes
- Context cues

This produces a feature representation that is small but semantically meaningful.

C. Text Generation

The text generator interprets the extracted features and formulates captions in a natural, sequential flow.

D. Streamlit-Based Web Interface

The system includes:

- Image upload option
- Preview display
- Caption output in real time.

This enhances ease of use and accessibility.

IX. SYSTEM ARCHITECTURE

The system architecture is designed for clarity, modularity and CPU-friendly execution. It separates concerns into frontend, preprocessing, visual feature extraction, text generation, and output/display modules so each part can be implemented, tested, and improved independently.

A. ARCHITECTURE OVERVIEW

The architecture consists of the following major layers:

1. User Interface (Frontend)
 - Streamlit-based web UI for image upload, preview, and caption display.
 - Accepts common image formats (JPG, PNG).
 - Shows progress and final caption to the user.
2. Preprocessing Module
 - Standardizes images: resizing to the model input resolution, normalization of pixel values, and conversion to the appropriate color channels.
 - Provides error handling for unsupported formats and basic image cleaning (crop/pad if necessary).
 - Exposes preprocessed tensors (or numpy arrays) to the encoder.
3. Visual Feature Extractor (Compact Encoder)
 - A lightweight CNN-based or pretrained feature extractor that produces a compact embedding vector representing objects, textures, spatial layout and global context.
 - Designed to run efficiently on CPU (reduced parameter count, smaller intermediate activations).
 - Optionally uses a small convolutional backbone with global pooling to create fixed-size feature vectors.
4. Text Generator (Sequence Module)
 - A compact sequential text generator that consumes the visual embedding and produces a caption token-by-token.

- Uses a small recurrent or transformer-lite block optimized for CPU inference (short beam/greedy decoding for speed).
- Includes vocabulary, tokenizer, and decoding strategy (greedy or beam-size 2–3).

5. Integration & Control Logic

- Glue layer between encoder and generator to pass features, manage decoding parameters, and handle exceptions.
- Performs any necessary reshaping, attention-compatible projections, or context concatenation.

6. Output Module

- Post-processes generated tokens into a readable sentence (detokenization, capitalization, punctuation).
- Displays the caption on the Streamlit UI and returns it via the API (if implemented).

7. Logging & Diagnostics

- Lightweight logging of input image IDs, timestamps, inference time per image, and errors.
- Optional local storage for anonymized diagnostic logs (not saving user images).

Non-functional considerations:

- Latency target: < 0.5 seconds per image on a typical laptop CPU (actual depends on model and hardware).
- Memory footprint: minimized by using compact models, smaller batch sizes and efficient data types.
- Privacy: images are processed locally (no upload to third-party servers) and not stored unless explicitly permitted.

B. OTHER FUNCTIONAL PARTS

- Streamlit Frontend- Makes it easier for users to interact Does away with complicated deployment steps Provides a real-time captioning feature
- Backend Engine - Takes care of image processing, feature extraction, running models, and generating captions. The backend guarantees quick execution and smooth integration.
- Compact Model Storage - The system works with adjusted model files that are good for academic, prototype, and demonstration use. This keeps memory use low while keeping caption accuracy high.

X. USE CASES

Automatic image captioning has become a big thing in modern computing, especially in areas where visual content needs text explanation. The uses of this project span different fields:

A. Support Technology for the Blind

The system can assist blind people in understanding visual content by changing images into spoken captions. Combining the model with text-to-speech tools can improve accessibility applications.

B. Social Media Automation

Content creators can get fast and suitable captions for images which lessens work and boosts interaction; the system can help platforms in auto-suggesting captions.

C. Digital Asset Management

Organizations may apply this solution for automatic tagging and categorization of images to enhance searchability and database organization.

D. E-commerce and Retail Platforms

Automatic captioning assists in describing products, particularly when visual details require textual representation for customers.

E. Surveillance and Security

Image captioning can help descriptive monitoring in surveillance systems by enabling textual logs of captured footage;

F. Education and Research

Students and researchers may use such models for understanding deep learning concepts conducting experiments or improving visual-language applications.

XI. PERFORMANCE EVALUATION

Caption quality was evaluated using a custom dataset and standard metrics:

- BLEU Score
- METEOR Score
- ROUGE-L Score

Results demonstrate strong semantic alignment and low caption latency despite CPU-only inference.

XII. COMPARATIVE ANALYSIS

The following table compares the proposed system with typical traditional captioning approaches.

Table I : Comparison of Image Captioning Approaches

Feature	Heavy Models	Transformer Models	Proposed System
Hardware	High-end GPU	Very High	CPU only
Speed	Moderate	Fast (GPU only)	Fast on CPU
Deployment	Complex	Very Complex	Simple
Power Use	High	High	Low
Practicality	Medium	High	Very High

This comparison highlights that the proposed system is ideal for academic, educational, and small-scale applications where efficiency and ease of use are more important than large-scale training performance.

XII. DISCUSSION

The study shows that it is possible to achieve meaningful caption generation without depending on large-scale architectures. The system balances both accuracy and efficiency, making it ideal for real-time academic demonstrations and lightweight applications..

XIV. LIMITATIONS

- Limited generalization on rare scenes
- Dependent on pretrained components
- English-only captions
- No automated learning from new data

XV. ETHICAL CONSIDERATIONS

- No storage of user images
- Neutral caption outputs
- Not intended for sensitive image analysis
- No harmful use cases supported

XVI. FUTURE SCOPE

- Multilingual captioning
- Larger datasets for training
- Real-time camera-based captioning

- Enhanced user interface
- Integration with mobile platforms

XVI. CONCLUSION

This work presents a practical and CPU-friendly image captioning system that generates contextually relevant descriptions with minimal computational overhead. The approach is suitable for small-scale deployment, academic usage, and accessibility tools. The results confirm that lightweight models can deliver effective performance without requiring high-end hardware.

XVII. REFERENCES

- Anderson, P., Fernando, B., Johnson, M., & Gould, S. (2018). Bottom-up and top-down attention for image captioning and visual question answering. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 6077–6086.
- Chen, X., Fang, H., Lin, T. Y., Vedantam, R., Gupta, S., Dollár, P., & Zitnick, C. L. (2015). Microsoft COCO captions: Data collection and evaluation server. *arXiv:1504.00325*.
- Donahue, J., Anne Hendricks, L., Guadarrama, S., Rohrbach, M., Venugopalan, S., & Saenko, K. (2015). Long-term recurrent convolutional networks for visual recognition and description. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2625–2634.
- Fang, H., Gupta, S., Iandola, F., Srivastava, R. K., Deng, L., & Zitnick, C. L. (2015). From captions to visual concepts and back. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1473–1482.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.
- Karpathy, A., & Fei-Fei, L. (2015). Deep visual-semantic alignments for generating image descriptions. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3128–3137.
- Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. *arXiv:1412.6980*.
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 1097–1105.
- Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., & Zitnick, C. (2014). Microsoft COCO: Common objects in context. *Proceedings of the European Conference on Computer Vision (ECCV)*, 740–755.
- Lu, J., Yang, J., Batra, D., & Parikh, D. (2017). Neural baby talk. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 7219–7228.
- Rennie, S. J., Marcheret, E., Mroueh, Y., Ross, J., & Goel, V. (2017). Self-critical sequence training for image captioning. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 7008–7024.
- Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., & Rabinovich, A. (2015). Going deeper with convolutions. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1–9.
- Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and tell: A neural image caption generator. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3156–3164.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., & Bengio, Y. (2016). Show, attend and tell: Neural image caption generation with visual attention. *Proceedings of the International Conference on Machine Learning (ICML)*, 2048–2057.
- Yao, T., Mei, T., & Rui, Y. (2016). Boosting image captioning with attributes. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 4894–4902.
- Zhang, Y., Hare, J., & Prügel-Bennett, A. (2018). Learning to count objects in natural images for image captioning. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Zhou, L., Kalantidis, Y., Chen, X., & Rohrbach, M. (2020). Unified vision-language pre-training for image captioning and VQA. *Proceedings of the European Conference on Computer Vision (ECCV)*, 465–481.
- Cornia, M., Baraldi, L., & Cucchiara, R. (2020). Meshed-memory transformer for image captioning. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 10578–10587.

-
- Huang, L., Wang, W., Chen, J., & Wei, X. (2019). Attention on attention for image captioning. Proceedings of the IEEE International Conference on Computer Vision (ICCV), 4634–4643.