

Visual Search Using Vision-Language Models (VLMs)

Ankit Kumar Gupta

Dept. of CSE

bwubta22544@brainwareuniversity.ac.in

Brainware University

Aditya Tiwari

Dept. of CSE

bwubta22184@brainwareuniversity.ac.in

Brainware University

Hirakjyoti Biswas

Dept. of CSE

bwubta22461@brainwareuniversity.ac.in

Brainware University

Abstract

Accurate image retrieval is a vital step in ensuring efficient digital experiences across various domains. Leveraging advancements in artificial intelligence, this study implements Vision-Language Models (VLMs) for automated visual search using multimodal learning. By processing annotated datasets through models like CLIP, PaLM-E, and ImageBind, the system achieves exceptional precision in text-to-image retrieval, similarity search, and zero-shot classification. Unlike conventional search methods that rely on metadata, our approach enables content-based retrieval, allowing users to find images based on semantic understanding. This paper also explores optimization strategies for real-time processing, making visual search scalable for large-scale applications. Furthermore, this research highlights the potential of VLMs in emerging areas like augmented reality, autonomous systems, and assistive technologies.

I. INTRODUCTION

Visual search has emerged as a transformative technology, enabling users to find images based on textual queries or reference images. Unlike traditional search techniques that rely on metadata and keyword tagging, modern approaches integrate computer vision and natural language processing (NLP) to understand semantic relationships between images and text. The rapid development of deep learning architectures, particularly Transformer-based models, has significantly enhanced the capabilities of VLMs, allowing them to generalize across diverse datasets with minimal supervision. This paper explores the recent advancements in multimodal learning and evaluates the effectiveness of VLMs in real-world applications. Additionally, we discuss the computational trade-offs of deploying these models at scale, particularly in edge computing environments, where resources are limited. By leveraging contrastive learning

techniques and fine-tuned embeddings, we propose a robust framework for efficient and accurate visual search.

Here are the roles for the three teammates:

- **Hirakjyoti Biswas:** Worked on **fine-tuning CLIP for text-to-image retrieval** and **implementing cosine similarity search for queries**.
- **Aditya Tiwari:** Focused on **data preprocessing & augmentation**, as well as **storing embeddings in FAISS for fast search**.
- **Ankit Kumar Gupta:** Handled **performance benchmarking across different models** and **optimization for real-time processing**.

II. METHODOLOGY

A. Approach

1. Data Acquisition: Use open-source image datasets (COCO, ImageNet).
2. Feature Extraction: Implement CLIP, ResNet, FAISS for embedding generation.
3. Query Processing: Convert input text/images into embeddings.
4. Similarity Computation: Perform cosine similarity search.
5. Model Fine-Tuning: Train models with additional domain-specific data for improved accuracy.
6. Performance Benchmarking: Compare retrieval accuracy and processing speed across different models.

B. Tools & Technologies Used

Deep Learning Models: CLIP, PaLM-E, ImageBind

Libraries: PyTorch, OpenCV, FAISS

Development Environment: Python, Jupyter Notebook

Hardware Requirements: High-performance GPUs for efficient training and inference.

III. PROJECT IMPLEMENTATION

A. Steps Taken

1. Data preprocessing & augmentation.
2. Fine-tuning CLIP for text-to-image retrieval.
3. Storing embeddings in FAISS for fast search.
4. Implementing cosine similarity search for queries.
5. Testing model accuracy with benchmark datasets.
6. Optimization for real-time processing.
7. Deployment of visual search system as an API for scalability.

B. Milestones Achieved

- [X] Developed CLIP-based text-to-image search.
- [X] Implemented image similarity search using FAISS.
- [X] Achieved zero-shot classification for unseen images.
- [X] Benchmarked model performance on real-world datasets.
- [X] Deployed a prototype system for user testing.

IV. RESULTS AND ANALYSIS

A. Findings & Observations

- CLIP's multimodal embeddings significantly improve text-based image search accuracy.
- FAISS indexing reduces retrieval time for large-scale datasets.
- Zero-shot classification achieves high accuracy in object recognition tasks.
- Fine-tuned models exhibit improved accuracy in domain-specific retrieval tasks.
- Real-time implementation enables interactive visual search with minimal latency.

V. CHALLENGES AND SOLUTIONS

A. Computational Limitations

- Challenge: High computational cost due to large model size.
- Solution: Used model quantization and optimized inference frameworks.

B. GPU not detected

- Challenge: Without the specific software installed the GPU was not detected and the model loaded on CPU.
- Solution: Installed Nvidia CUDA 12.1 which is compatible with the version of python libraries which we were using.

C. Dataset

- Challenge: Normally the most popular datasets are very large and too big to install on device.
- Solution: We searched for a smaller dataset which covers a wide variety of content.

D. Scalability Issues

- Challenge: Large-scale image retrieval requires extensive storage and indexing.
- Solution: Deployed FAISS with efficient hierarchical indexing for fast retrieval.

VI. CONCLUSION AND FUTURE WORK

A. Conclusion

This study successfully implemented a multimodal visual search system using state-of-the-art Vision-Language Models. The combination of CLIP, FAISS, and AI-driven embedding search significantly improved image retrieval accuracy.

B. Future Enhancements

1. Real-time search optimization for low-latency applications.
2. Expanding dataset diversity for better generalization.
3. Integrating multi-modal retrieval (e.g., combining images, text, and video search).
4. Developing mobile-friendly visual search applications.
5. Enhancing robustness to adversarial queries.

REFERENCES

- [1] A. Radford et al., 'Learning Transferable Visual Models From Natural Language Supervision,' OpenAI, 2021.
- [2] J. Johnson et al., 'Billion-scale similarity search with GPUs,' Facebook AI Research, 2017.
- [3] X. Wang et al., 'Multimodal Learning for Image-Text Retrieval,' IEEE Transactions on Neural Networks, 2022.

[4] Y. LeCun et al., 'Deep Learning for Visual Recognition,' Nature, 2021.