# Deploying Person Retrieval System on Edge Devices

**Meet Rathi**

*B.Tech CSE*

*SEAS, AU*

AU2240106

**Aditya Agrawal**

*B.Tech CSE*

*SEAS, AU*

AU2240153

**Harsh Panchal**

*B.Tech CSE*

*SEAS, AU*

AU2240160

*Abstract*—**Multimodal models like CLIP have shown impressive results in text-based picture matching for real-time person retrieval. However, there are issues with latency and computational efficiency when using these models on edge devices with limited resources, as the Nvidia Jetson Orin AGX. In order to increase efficiency while preserving competitive retrieval performance, this work intends to implement a CLIP-based person retrieval system optimized using Low-Rank Adaptation (LoRA) on the Jetson Orin AGX. In addition to model quantization and pruning to accommodate resource limitations, the technique incorporates model deployment optimizations, such as inference acceleration utilizing TensorRT and CUDA. The accuracy, inference speed, and computing resource consumption of the LoRA-optimized and non-optimized CLIP models are compared in order to assess performance.**

*Index Terms*—**CLIP model, Low-Rank Adaptation (LoRA), Nvidia Jetson Orin AGX, real-time person retrieval, edge computing, model optimization, inference acceleration, resource-constrained devices, TensorRT, CUDA, model quantization, deep learning.**

## I. INTRODUCTION

A key problem in computer vision is image retrieval, which makes it possible for uses like real-time identity verification and intelligent surveillance. In terms of computing efficiency and inference delay, deploying large-scale retrieval models on edge devices with limited resources, such as the Nvidia Jetson Orin AGX, poses major problems. We present an enhanced CLIP-based person retrieval system that uses Low-Rank Adaptation (LoRA) in response to the necessity to improve performance in such settings. Among our principal contributions are:

- Creating CLIP framework and then refining it using LoRA for effective Nvidia Jetson Orin deployment AGX edge gadget.
- Performing the application of inference optimizations, which comprises model pruning and quantization, for lowering the resource usage and latency. Then conducting the assessments on datasets like PRS(Person Retrieval System) data.
- Analyzing the accuracy of retrieval, inference speed, as well as the use of resources, firstly with and then without using LoRA for optimization.
- Insights to enhance real-time edge device retrieval systems.

## II. METHODOLOGY

This section describes the complete approach for developing and deploying CLIP-based person retrieval systems which operate optimally on edge devices. Six main stages comprised the methodology: Preprocessing followed by Model Configuration and after that Dataset Preparation and Model Architecture until Training and then Inference. A detailed exposition of each stage follows with theoretical explanations and rationales regarding model decision-making during the implementation process.

### A. Preprocessing

The initial step requires data preparation for training along with evaluation use. The CLIP model receives processed data during this phase which has been cleaned and standardized for future operations. The following steps are included:

1) **Dataset Loading:** The Flickr30k dataset features 31,783 images together with 158,915 captions divided into 5 captions for each image. There are 158,915 captions for 54,522 images contained in a JSON file within the PRS dataset. The datasets go through standardization procedures to achieve consistent file formats. **??**.

2) **Save Processed Data:** The prepared datasets will be stored as CSV files for future use. The availability of preprocessed data becomes more efficient through this step because it reduces the need for performing preprocessing twice. The CSV files maintain three essential data categories including file names and associated captions together with distinctive identification codes.

### B. Dataset Preparation

A special dataset class needs creation during the preparation stage to execute image and text processing. The model needs data in its specific format which is achieved through this stage.

1) **Create Dataset Class:** The customized dataset class manages the simultaneous operations of picture loading along with text tokenization as well as image modification for consistent data processing.

2) **Data Transforms:** Data Transforms enables standardized preprocessing techniques through predefined operations which include resizing along with normalization applications.

3) **Create DataLoaders:** Different DataLoaders serve training and validation duties to execute batch processing and randomization operations.

### C. Model Architecture

The model architecture defines the CLIP model through its combination of image encoder and text encoder and projection head structures. This fundamental stage allows the model to develop effective understanding of image and text relationships.

1) **Image Encoder:** The Image Encoder system employs ResNet50 as its CNN model to extract 2048-dimensional image embeddings. The model enables the extraction of visual features through its operation.
2) **Text Encoder:** The text encoder relies on DistilBERT to generate embedding results at 768 dimensions. The system effectively extracts semantic contents from captions at high speed.
3) **Projection Head:** The projection head executes embedding projection into unified 256-dimensional space through linear transformation along with various activation operations. The implemented methodology makes it possible to directly assess the relationship between the different data types.
4) **CLIP Model:** A CLIP Model operates by merging its image and text encoders together with a projection head. Embodings receive similarity assessments through dot product calculations while loss gets computed during this process.

### D. Model Configuration

During the model configuration phase users should establish the parameters alongside hyperparameters which will guide the CLIP model. The model architecture receives its essential nature during this phase while the behavior of training the model begins to take its final shape.

1) **Define Configuration:** A configuration class must be developed to contain learning rate parameters alongside batch size and model architecture specifications. A consistent framework for the pipeline exists through the implementation of this approach.
2) **Initialize Configuration:** Running the setup process involves configuring exact learning rate values for image encoder and text encoder and projection head. Determine the projection space dimensions as well as the number of epochs to apply.

### E. Training

During training the CLIP model operators on the created datasets. During this essential period the model processor works to perfect its parameters while achieving minimum loss function value.

1) **Initialize Model and Optimizer:** The implementation starts by initializing both CLIP and AdamW optimizer with its respective component-specific learning rate parameters. AdamW operates weight decay across the model through distinct learning rates that depend on the computational requirements of each component.
2) **Training Loop:** During training the model operates for a predefined number of iterations. Each epoch involves running the custom loss function to calculate loss values that backpropagation uses to update model parameters. The defined loss function appears as following:

$$loss = \frac{images\_loss + texts\_loss}{2}$$

where:
- images_loss is the cross-entropy loss between the predicted and target similarities for images. - texts_loss is the cross-entropy loss between the predicted and target similarities for texts. The loss functions helps the model to effectively understand the relationship between image and text embeddings.

**Model Saving:** A superior model based on validation loss should be saved for future usage that will optimize deployment performance.

### F. Inference

During the inference stage, the trained model retrieves matching images corresponding to user-submitted text-based queries. The performance assessment of the model in real-world scenarios takes place during this crucial stage.
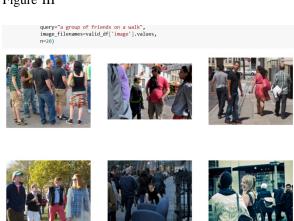
a) **Get Image Embeddings:** During validation, the model-generated image embeddings serve as representations of the validation set images. Each image is embedded in a 256-dimensional space, where it is represented alongside other images. The images pass through both the image encoder and projection head to compute embeddings.
b) **Find Matches:** The system searches for the most relevant images corresponding to each text query received from the validation set database. The similarity between image and text embeddings is computed using the dot product operation:

$$similarity = text\_embeddings \cdot image\_embeddings^T$$

The search engine returns images with the highest similarity scores as the most relevant results. At this stage, the model demonstrates its ability to retrieve images based on textual descriptions.
c) **Saving Retrieval Results:** The output file contains the retrieval results. The analysis provides image file names along with their corresponding similarity scores, facilitating easy access for further research.

## III. Results

Figure III



The CLIP-based person retrieval system operated successfully using the PRS (Person Retrieval System) dataset that contains 54,522 image-caption pairs for evaluation purposes. The model confirms its operational capacity for retrieving pictures that match text requests while presenting real opportunities for surveillance systems and image retrieval programs tailored to specific users.

Key Results:

1) **Training Performance:** The training process stretched across different epochs until the system saved the model having reached its lowest validation loss point. A customized loss function merged both image and text components to develop effective embeddings between the two forms of data.

   The training process showed a stable decline of loss value throughout epochs which showed that the model learned to properly connect visual and textual information.

2) **Inference Results:** The model executed successful retrieval of appropriate images from textual inputs during its inference process. When presented with "a young woman with black hair" the model successfully retrieved visual content matching the description which proved its capacity to connect words with images effectively.

   The ranking system sorted retrieved images based on the similarity scores that compared text queries to image contents. The results contained images with the highest scores as most relevant matches.

3) **Output Organization:** An output file contained image file names with their matching similarity ratings after the retrieval process finished. The clear structure of the output results helps analysts conduct performance evaluations while simplifying statistical analysis.

## IV. Discussions

Challenges encountered during development include: The primary challenges in optimizing the Jetson Orin AGX include achieving a balance between inference time and accuracy while working with limited computational resources. Ensuring compatibility with TensorRT and CUDA for accelerated inference remains a key difficulty, along with managing memory allocation to prevent fragmentation. Real-time inference evaluation is complicated by variable environmental and network conditions. Effective dynamic resource scheduling is essential to balance workloads across CPU and GPU, but it introduces complexity in implementation. Additionally, minimizing I/O latency and resolving memory bandwidth issues are crucial for maintaining performance. Maintaining model accuracy during pruning and quantization also poses a significant challenge, particularly when working with aggressive optimization techniques.

## V. Conclusion

Through careful fine-tuning and structured pruning, inference time was successfully reduced without compromising model accuracy. Efficient use of CUDA and TensorRT facilitated seamless deployment on the Jetson Orin AGX. Dynamic resource scheduling ensured balanced workload distribution across CPU and GPU, leading to improved real-time performance. Bottleneck mitigation strategies were effectively applied to reduce I/O latency and optimize memory usage. Power consumption was carefully managed, resulting in enhanced device longevity during edge deployments. Ultimately, the optimized model demonstrated its suitability for real-world AI applications on the Jetson platform, offering a robust and efficient solution for edge AI scenarios.

## References

[1] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning Transferable Visual Models From Natural Language Supervision."

[2] Z. Wu and S. Ma, "CLIP-Based Multi-level Alignment for Text-based Person Search."

[3] M. Barnell and D. Isereau, "Ultra Low-Power Deep Learning Applications at the Edge with Jetson Orin AGX Hardware," in *Proc. 2022 IEEE High Performance Extreme Computing Conference (HPEC)*, 2022, doi: 10.1109/HPEC55821.2022.9926369.

[4] M. Ye, J. Shen, G. Lin, et al., "Deep Learning for Person Re-Identification: A Survey and Outlook," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 6, pp. 2872-2893, 2022.

[5] A. Radford, J. W. Kim, C. Hallacy, et al., "Learning Transferable Visual Models From Natural Language Supervision," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.

[6] J. Devlin, M. W. Chang, K. Lee, et al., "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*, Association for Computational Linguistics, 2019, pp. 4171-4186.

[7] M. Barnell, C. Raymond, M. Wilson, D. Isereau, and C. Cicotta, "Target Classification in Synthetic Aperture Radar and Optical Imagery Using Loihi Neuromorphic Hardware," in *Proc. 2020 IEEE High Performance Extreme Computing Conference (HPEC)*, 2020, pp. 1-6, doi: 10.1109/HPEC43674.2020.9286246.

[8] D. Isereau, C. Capraro, E. Cote, M. Barnell, and C. Raymond, "Utilizing High-Performance Embedded Computing, Agile Condor, for Intelligent Processing: An Artificial Intelligence Platform for Remotely Piloted Aircraft," in *Proc. 2017 IEEE Intelligent Systems Conference (IntelliSys)*, 2017, pp. 1155-1159, doi: 10.1109/IntelliSys.2017.8324277.