

AI-Driven Image Generation and Virtual Try-On for Personalized Fashion Experiences

Ashwini K. Sukanawar

M. Tech Student of Department of Computer Science & Engineering, Shri Balasaheb Mane Shikshan Prasarak Mandal's, Ashokrao Mane Group of Institutions, Vathar, Kolhapur, India

Abstract - The integration of Retrieval-Augmented Generation (RAG) with Large Language Models (LLMs) presents a transformative approach to generating personalized and contextually relevant images that cater to specific user preferences. This project aims to harness the synergistic potential of RAG and LLMs to develop a robust and scalable image generation pipeline that seamlessly blends state-of-the-art natural language processing with advanced computer vision techniques. The process begins by utilizing a RAG model, which combines the strengths of retrieval-based methods and generative models to produce high-quality images that are not only coherent with the input prompts but also enriched with context from external knowledge sources. Following the image generation, a dedicated preprocessing module is employed to resize and optimize the images, ensuring they meet the quality standards required for subsequent integration. The next critical phase involves the detection of human upper bodies in photographs using Haar Cascade classifiers, a machine learning-based approach known for its efficiency in real-time object detection. The accurate identification of the upper body regions is crucial for the next step, where the generated images are overlaid onto these detected regions using OpenCV, a powerful computer vision library. This integration ensures that the images are aligned precisely with the contours of the human body, creating a visually realistic and aesthetically pleasing effect. To facilitate user interaction and deployment, the entire process is encapsulated within a Flask framework, which serves as the backbone of the application's architecture. The Flask framework not only handles the backend processing, including API requests and image processing tasks, but also supports a user-friendly frontend interface, allowing users to interact with the system effortlessly.

Keywords: Augmented Generation (RAG), Large Language Models (LLMs), image generation, Machine Learning, Open CV, Deep Learning.

I. INTRODUCTION

The rapid advancements in AI and deep learning have opened up new possibilities for image generation and computer vision. The project aims to combine the strengths of Retrieval-Augmented Generation (RAG) and Large Language Models (LLMs) with computer vision techniques to create a system that can generate and overlay images onto human bodies in a realistic and user-friendly manner. This innovative approach not only enhances the quality of generated images but also provides practical applications in various fields, such as virtual try-ons and personalized fashion recommendations. By utilizing a Flask web framework, the system is designed to be accessible and scalable, catering to a wide range of users.

The convergence of RAG and LLM technologies with computer vision marks a significant step forward in the realm of personalized image generation and augmentation. Traditional image generation methods, while effective, often struggle to incorporate context-sensitive details that align with user-specific requirements. By integrating RAG, which combines the retrieval of relevant data from vast repositories with the generative capabilities of LLMs, this project is able to create images that are not only visually appealing but also contextually appropriate and highly personalized. Moreover, the incorporation of appealing but also contextually appropriate and highly personalized.

The incorporation of computer vision techniques, particularly through the use of Haar Cascade classifiers and OpenCV, allows for the dynamic adaptation of these generated images to fit the contours and dimensions of human bodies in photographs. This approach addresses the need for a more immersive and interactive user experience, enabling applications that go beyond static image generation to offer real-time, adaptable visual content. The Flask framework ensures that these complex processes are packaged into a user-friendly application, making advanced AI-driven image generation accessible to users across different technical proficiencies and industries.

II. EXISTING SYSTEM

1) Standalone Image Generation:

Description: The majority of existing image generation systems focus on creating images from text inputs or prompts. These systems typically use Generative Adversarial Networks (GANs) or other deep learning models to generate images that match the given description.

Limitations: While these systems are capable of producing high-quality images, they are primarily designed to generate standalone images without any context. The generated images exist in isolation and are not intended to be integrated with real-world elements, such as human figures or environments.

2) Contextual Image Generation:

Description: Some advanced systems attempt to generate images that consider contextual information, such as style, theme, or specific object placement. These systems may use techniques like style transfer or conditional GANs to create images that align with certain contextual cues.

Limitations: Although these systems are more advanced than standalone generators, they often struggle to adapt the generated images to specific real-world contexts, such as fitting an image onto a human body or aligning it with complex backgrounds. They still largely produce images that are contextually relevant but not interactively adaptable.

3) Manual Integration with Real-World Elements:

Description: Some existing systems allow for the manual overlay of generated images onto human figures or other real-world elements. This process often involves using image editing software or basic computer vision tools to manually place and adjust the generated images onto a given background or figure.

Limitations: The manual nature of this process is time-consuming and requires significant human intervention. It lacks automation and scalability, making it impractical for large-scale applications. Moreover, the accuracy and realism of the overlay depend heavily on the user's skill, often resulting in suboptimal outcomes.

4) Simplistic Computer Vision Integration:

Description: A few systems attempt to integrate generated images with human figures using basic computer vision techniques. For instance, simple bounding box methods or keypoint detection may be used to roughly align generated images with body parts.

Limitations: These systems tend to use overly simplistic approaches that do not fully leverage the capabilities of advanced AI or deep learning. As a result, the alignment of images is often imprecise, leading to unnatural or unrealistic overlays that detract from the user experience.

5) Limited Use of AI for Contextual Understanding:

Description: Existing systems that use AI for contextual understanding are often limited in their ability to comprehend complex real-world scenarios. For example, they may generate an image of clothing based on a text prompt but lack the capability to accurately position that clothing on a human figure in a photograph.

Limitations: The AI models used in these systems may lack the sophistication required to fully understand and interpret the context in which the generated images will be used. This results in outputs that may be contextually relevant in theory but fall short in practical application.

III. PROPOSED METHODOLOGY

1) Image Generation Using RAG and LLM:

Objective: To generate high-quality and contextually relevant images based on user inputs or prompts.

Approach:

- **API Integration:** Begin by generating an API key for OpenAI or another LLM model provider. This API will be used to access the model for image generation.
- **RAG Model Implementation:** Implement a Retrieval-Augmented Generation model that combines retrieval-based methods with generative capabilities. The model will retrieve relevant context from a large corpus of data and use this information to enhance the image generation process.
- **Text-to-Image Generation:** Utilize the LLM to interpret the user's prompts or text inputs and generate corresponding images. The model will be fine-tuned to ensure that the generated images are not only visually appealing but also align with the user's specific requests and contextual needs.

2) Image Preprocessing and Optimization:

Objective: To ensure that the generated images meet the necessary quality standards and are suitable for overlaying onto human figures.

Approach:

- **Resizing and Cropping:** Develop an image processing model that preprocesses the generated images. This includes resizing and cropping the images to fit the dimensions required for overlaying onto detected human bodies.
- **Quality Enhancement:** Implement techniques to enhance image quality, such as sharpening, contrast adjustment, and noise reduction, ensuring that the images remain clear and visually consistent.
- **Format Conversion:** Convert the images into the appropriate format and resolution needed for further processing in the computer vision pipeline.

3) Human Body Detection Using Haar Cascade Classifier:

Objective: To accurately detect the upper body of humans in photographs, which is essential for correctly overlaying the generated images.

Approach:

- **Dataset Collection:** Compile a diverse dataset of human images featuring various poses, body types, and backgrounds to train the Haar Cascade classifier.
- **Classifier Training:** Train the Haar Cascade classifier to detect the upper body of humans in different images. The classifier will learn to identify key features such as the head, shoulders, and torso.
- **Performance Optimization:** Optimize the classifier by adjusting parameters such as the scale factor, minNeighbors, and minSize to improve detection accuracy and reduce false positives.

4) Image Overlay Using OpenCV:

Objective: To integrate the processed images with the detected human figures, ensuring accurate and aesthetically pleasing overlays.

Approach:

- **Integration with Haar Cascade:** Use OpenCV to integrate the output of the Haar Cascade classifier with the generated images. This involves aligning the images with the detected body parts in a way that maintains realism and natural appearance.
- **Image Transformation:** Apply geometric transformations such as rotation, scaling, and translation to the images to match the orientation and size of the detected body regions.
- **Overlay Logic:** Develop the logic to handle the overlay process, ensuring that the images are correctly positioned

and blended with the background, maintaining a natural look.

5) System Integration with Flask Framework:

Objective: To create a user-friendly web application that allows users to interact with the system, generate images, and view the final output.

Approach:

- **Backend Development:** Develop the backend of the application using Flask, a lightweight and flexible web framework. The backend will handle API requests, manage the image generation and processing pipeline, and communicate with the frontend.
- **Frontend Development:** Design a user-friendly frontend that allows users to input prompts, upload photos, and view the final output. The frontend will be integrated with the Flask backend to provide a seamless user experience.
- **Real-Time Processing:** Ensure that the system is capable of processing images in real-time, providing quick and responsive feedback to users. This includes optimizing the pipeline to handle concurrent requests without significant delays.

6) Evaluation and Testing:

Objective: To evaluate the performance of the entire system and ensure it meets the required standards for accuracy, efficiency, and user satisfaction.

Approach:

- **Model Testing:** Test the RAG, LLM, and Haar Cascade models individually and in combination to ensure they perform well in generating and overlaying images.
- **User Testing:** Conduct user testing sessions to gather feedback on the system's usability, responsiveness, and overall quality of the generated outputs.
- **Performance Metrics:** Measure the system's performance using relevant metrics such as accuracy, precision, recall, response time, and user satisfaction scores. Use this data to identify areas for improvement and further optimization.

IV. CHALLENGES AND LIMITATIONS

1) Complexity of Multi-Model Integration:

- **Challenge:** Integrating multiple advanced models such as RAG, LLMs, and Haar Cascade classifiers into a cohesive pipeline can be complex.

- **Impact:** Ensuring seamless communication between these components requires careful design and debugging, which can be time-consuming and prone to errors.

2) Accuracy of Human Body Detection:

- **Challenge:** Achieving high accuracy in detecting the upper body of humans in various images using Haar Cascade classifiers.
- **Impact:** Inaccurate detection can lead to misaligned overlays, reducing the realism and effectiveness of the final output.

3) Quality of Generated Images:

- **Challenge:** Maintaining high quality and contextual relevance in the images generated by the RAG and LLM models.
- **Impact:** Low-quality or irrelevant images can diminish the user experience and reduce the utility of the system.

4) Real-Time Processing Demands:

- **Challenge:** Ensuring that the entire pipeline, from image generation to overlay, operates efficiently in real-time.
- **Impact:** High processing demands may lead to delays, affecting the responsiveness of the system and potentially causing user frustration.

5) Dataset Diversity for Training:

- **Challenge:** Training the Haar Cascade classifier on a diverse dataset to ensure it can handle various human body types, poses, and environmental conditions.
- **Impact:** A lack of diversity in the training data can result in poor generalization, leading to inaccuracies in body detection in real-world scenarios.

6) Balancing Model Performance with Resource Usage:

- **Challenge:** Balancing the need for high model performance with the constraints on computational resources, particularly in a web-based application.
- **Impact:** High computational requirements may limit the scalability of the system and increase the cost of deployment.

Limitations:

1) Dependence on Pre-Trained Models:

- **Limitation:** The project relies heavily on pre-trained RAG and LLM models, which may not be perfectly suited to the specific context of image generation for body overlays.

- **Impact:** The pre-trained models might require significant fine-tuning to achieve optimal results, which could be resource-intensive.

2) Limited Flexibility in Image Types:

- **Limitation:** The system may be limited in the types of images it can generate and overlay, particularly if the images are highly specific or complex.
- **Impact:** This could restrict the applicability of the system to certain use cases, limiting its versatility.

3) Scalability Constraints:

- **Limitation:** The computational demands of real-time image generation, processing, and overlay may pose challenges in scaling the system to support a large number of users simultaneously.
- **Impact:** Scalability issues could hinder the system's adoption in high-traffic environments, such as popular e-commerce platforms or social media applications.

4) User Dependency for Input Quality:

- **Limitation:** The quality of the generated output heavily depends on the quality and specificity of the user inputs or prompts.
- **Impact:** Ambiguous or low-quality inputs may lead to suboptimal results, requiring users to have a certain level of expertise or understanding to achieve the best outcomes.

5) Challenges in Cross-Platform Integration:

- **Limitation:** Integrating the system with various platforms (e.g., mobile devices, different web browsers) can introduce compatibility issues.
- **Impact:** Ensuring consistent performance and appearance across different platforms may require additional development and testing efforts.

6) Ethical and Privacy Considerations:

- **Limitation:** The use of personal images for overlaying generated content raises ethical and privacy concerns, especially regarding the handling and storage of user data.
- **Impact:** Addressing these concerns requires robust data protection measures and clear user consent protocols, adding complexity to the system's design and implementation.

V. RESULTS AND DISCUSSION

1) Significance of the Project:

- **Advancement in AI Integration:** The project represents a significant leap in combining RAG (Retrieval-Augmented Generation) with LLMs (Large Language Models) for image generation. This integration allows for the creation of personalized, contextually relevant images that are tailored to user preferences.
- **Real-World Application:** Unlike traditional image generation systems that produce standalone images, this project focuses on generating images that can be seamlessly integrated into real-world contexts, specifically by overlaying them onto human figures. This opens up new possibilities in industries like fashion, advertising, and interactive media.

2) Challenges in Image Alignment:

- **Accurate Overlay:** One of the key challenges is ensuring that the generated images are accurately aligned with detected human figures. Misalignment can result in unrealistic or distorted images, which would undermine the effectiveness of the system.
- **Use of Haar Cascade Classifiers:** Haar Cascade classifiers are employed to detect the upper body of humans in images. The accuracy of this detection is critical, as it directly affects the quality and realism of the final output. The classifier must be well-trained to handle various scenarios and ensure precise detection.

3) Leveraging RAG and LLMs:

- **RAG Model:** The RAG model combines retrieval-based methods with generative capabilities, allowing it to produce images that are enriched with contextual information. This ensures that the generated images are not only high in quality but also relevant to the user's input.
- **Personalization through LLMs:** By incorporating LLMs, the system can generate images that are tailored to specific prompts or user needs, enhancing the personalization aspect of the image generation process.

4) Complexity of Technology Integration:

- **Multi-Technology Integration:** The project involves the integration of multiple advanced technologies, each with its own complexities. Balancing these technologies to work seamlessly together is a significant challenge.
- **Tuning the RAG Model:** The RAG model must be carefully tuned to ensure it generates images that are contextually appropriate and of high quality. This

requires a deep understanding of both the model and the data it operates on.

- **Training the Haar Cascade Classifier:** The Haar Cascade classifier requires extensive training on a diverse dataset to accurately detect human upper bodies. This is crucial for ensuring that the generated images align correctly with the detected figures.

5) Image Preprocessing:

- **Quality Maintenance:** The preprocessing step, which includes resizing and adjusting the images, is essential to maintain image quality. Improper preprocessing could lead to loss of detail or distortions, negatively impacting the final output.
- **Efficiency in Processing:** The preprocessing must be efficient to handle real-time demands, especially when integrated into a web application where response time is critical.

6) Flask Framework for Deployment:

- **Lightweight and Flexible:** Flask is chosen for its lightweight and flexible nature, making it suitable for developing the web application that will host the image generation system.
- **Integration Challenges:** Integrating such a complex pipeline within Flask involves managing API requests, real-time image processing, and ensuring that the system remains responsive and user-friendly.
- **User Interaction:** Flask provides the necessary tools to create an interactive and accessible frontend, allowing users to easily engage with the system and view the results of the image generation process.

7) Potential Impact and Applications:

- **Fashion Industry:** The project could revolutionize the fashion industry by enabling virtual try-ons and personalized fashion recommendations. Users could upload their photos and see how different outfits look on them in real-time.
- **Personalized Advertising:** Advertisers could use this technology to create highly personalized and interactive ad campaigns, where the generated images are tailored to individual users and overlaid onto their figures.
- **Interactive Media and Content Creation:** The system could be used in various interactive media applications, allowing users to generate and manipulate images that fit within specific contexts, enhancing the overall user experience.

VI. CONCLUSIONS

The integration of Retrieval-Augmented Generation (RAG) with Large Language Models (LLMs) and advanced computer vision techniques in this project marks a significant advancement in the field of AI-driven image generation and manipulation. By successfully combining these technologies, the project has demonstrated the potential to generate contextually relevant and personalized images that can be seamlessly overlaid onto human figures. This approach not only enhances the realism and utility of the generated images but also opens up new possibilities for applications in fashion, advertising, and interactive media, where personalized and context-aware content is increasingly in demand.

The development process highlighted several challenges, including the complexity of integrating multiple models, ensuring accurate human body detection, and maintaining high-quality image generation. However, through careful design, rigorous testing, and the use of robust frameworks like Flask, these challenges were effectively addressed. The result is a system that is both powerful and user-friendly, capable of delivering high-quality outputs in real-time. The project also underscores the importance of continuous model evaluation and refinement, as well as the need for diverse training datasets to ensure the system's adaptability to various real-world scenarios.

REFERENCES

1) RAG and LLMs:

- [1] Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., ... & Riedel, S. (2020). "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks." Proceedings of the 34th Conference on Neural Information Processing Systems (NeurIPS), 1-16.
- [2] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P.,... & Amodei, D. (2020). "Language Models are Few-Shot Learners." Advances in Neural Information Processing Systems (NeurIPS), 33, 1877-1901.
- [3] Guu, K., Lee, K., Tung, Z., Pasupat, P., & Chang, M. (2020). "REALM: Retrieval-Augmented Language Model Pre-Training." Proceedings of the 37th International Conference on Machine Learning (ICML), 8877-8888.
- [4] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2020). "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer." Journal of Machine Learning Research, 21(140), 1-67.
- [5] Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., ... & Yih, W. T. (2020). "Dense Passage

Retrieval for Open-Domain Question Answering." Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), 6769-6781.

2) OpenCV and Computer Vision:

- [6] Bradski, G. (2000). "The OpenCV Library." Dr. Dobb's Journal of Software Tools, 25(11), 120-125.
- [7] Kaehler, A., & Bradski, G. (2016). "Learning OpenCV 3: Computer Vision in C++ with the OpenCV Library." O'Reilly Media.
- [8] Dalal, N., & Triggs, B. (2005). "Histograms of Oriented Gradients for Human Detection." Proceedings of the 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 886-893.
- [9] Viola, P., & Jones, M. (2001). "Rapid Object Detection Using a Boosted Cascade of Simple Features." Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 511-518.

3) General AI and Deep Learning:

- [10] Goodfellow, I., Bengio, Y., & Courville, A. (2016). "Deep Learning." MIT Press.
- [11] He, K., Zhang, X., Ren, S., & Sun, J. (2016). "Deep Residual Learning for Image Recognition." Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 770-778.

AUTHOR'S BIOGRAPHY



Ms. Ashwini K. Suganawar,
M. Tech Student of Department of Computer Science & Engineering, ShriBalasaheb Mane Shikshan Prasarak Mandal's, Ashokrao Mane Group of Institutions (AMGOI), Vathar, Kolhapur, India.

Citation of this Article:

Ashwini K. Sukanawar, (2024). AI-Driven Image Generation and Virtual Try-On for Personalized Fashion Experiences. *International Research Journal of Innovations in Engineering and Technology - IRJIET*, 8(9), 112-118. Article DOI <https://doi.org/10.47001/IRJIET/2024.809014>
