
COSE474-2024F: Final Project Proposal

A text-based image editing tool combining CLIP, DETR and Stable Diffusion

2020320041 Seokmin Kim

1. Introduction

1.1. Motivation

Image editing has become a significant research topic in computer vision and artificial intelligence. Although recent advances in diffusion models have dramatically improved image generation capabilities, precisely editing specific objects within an image remains a challenging task. Traditional image editing methods either require extensive technical expertise or are limited to global modifications that affect the entire image.

This research proposes a novel framework that combines DETR's (Carion et al., 2020) object detection capabilities, CLIP's (Radford et al., 2021) text-image understanding, and Stable Diffusion's (Rombach et al., 2022) image editing power. Through this integration, our aim is to provide an intuitive and effective image editing method in which users can precisely modify specified objects through natural language instructions.

1.2. Problem definition

The primary challenge in text-guided object-specific image editing lies in three key aspects. First, accurately detecting and localizing the target object specified by the user's text prompt while maintaining awareness of the surrounding context. Second, generating a precise mask that isolates only the target object for editing. Third, applying the desired modifications to the masked region while preserving the visual coherence with unedited portions of the image.

1.3. contribution

This research introduces a novel framework for precise text-guided object-specific image editing that seamlessly integrates DETR, CLIP, and Stable Diffusion. The proposed approach develops an efficient object masking mechanism that combines DETR's object detection capabilities with CLIP's semantic understanding to accurately isolate target objects. Through a comprehensive evaluation using the COCO (Lin et al., 2014) dataset, the effectiveness of this approach is demonstrated using both quantitative metrics (FID, SSIM) and qualitative results. Furthermore, this framework provides an intuitive interface that enables users to edit spe-

Algorithm 1 Text-Guided Object-Specific Image Editing

Input: Image I , prompt P

Output: edited Image E

repeat

Object Detection

 Initialize $boxes, scores = DETR(I)$

 Initialize $relevantBoxes = filterBoxes(scores > threshold)$

Object-Text Matching

 Initialize $matchScores = []$

for box **in** $relevantBoxes$ **do**

$crop = extractCrop(I, box)$

$similarity = CLIPSimilarity(crop, P)$

$matchScores.append(similarity)$

end for

$bestBox = relevantBoxes[argmax(matchScores)]$

Mask Generation and Editing

$mask = generateMask(bestBox)$

$E = StableDiffusionInpaint(I, mask, P)$

until edit complete

return E

cific objects through natural language instructions while maintaining image coherence.

2. Methods

I propose the method addresses the challenges of text-guided object-specific image editing through a three-stage architecture integrating object detection, semantic understanding, and localized image editing. This approach is distinguished by its ability to accurately identify and edit specific objects while maintaining global image coherence. The key innovation lies in the synergistic combination of DETR's precise object detection, CLIP's semantic matching, and Stable Diffusion's controlled editing capabilities.

The framework begins with the Object Detection and Localization stage, which employs a DETR-ResNet50 backbone for initial object detection. The input images are processed through a custom transform pipeline for 800x800 resolution, where the module generates bounding boxes with confidence scores for detected objects.

In the Semantic Understanding and Matching stage, the sys-

tem utilizes CLIP’s vision and text encoders for semantic alignment. This stage computes similarity scores between detected objects and text prompts using cosine similarity matching for accurate object-text pairing. The implementation includes batch processing capabilities for efficient handling of multiple objects, enhancing the system’s overall performance.

The final stage, Localized Image Editing, creates precise binary masks from DETR’s bounding boxes and scales them to 512x512 resolution for compatibility with Stable Diffusion. The system employs the Stable Diffusion Inpainting Pipeline for targeted modifications while maintaining contextual consistency through controlled diffusion steps. This stage features adaptive object selection through the integration of CLIP similarity scores with DETR confidence scores, along with dynamic threshold adjustment for optimal object selection.

The mask generation process incorporates a coordinate transformation system for accurate alignment and boundary refinement for smooth mask edges. The controlled diffusion process includes custom scheduling for inpainting steps and context-aware guidance scale adjustment, ensuring high-quality results while preserving the original image’s integrity.

This comprehensive architecture effectively addresses the main challenges of object-specific editing by ensuring accurate object identification through multi-modal verification, maintaining precise spatial localization of edits, and preserving global image coherence while applying local modifications.

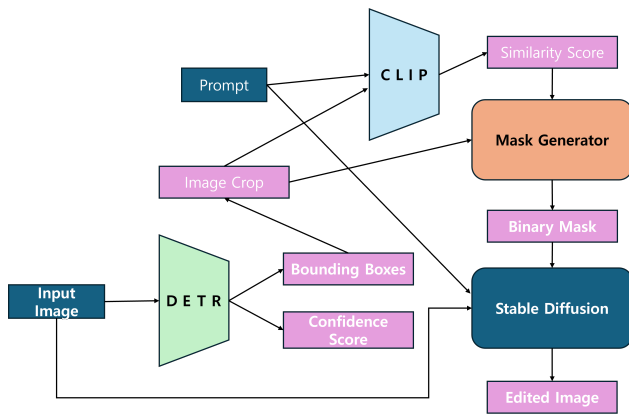


Figure 1. Architecture of the Text-Guided Object-Specific Image Editing Model

3. Experiments

3.1. Dataset

I choose the COCO 2017 validation dataset for the experiments, which provides a diverse collection of real-world images with comprehensive object annotations. The validation set contains approximately 5,000 images across 80 object categories, making it particularly suitable for assessing object-specific editing capabilities. The experiments focused on natural objects and common items to evaluate the model’s performance across varying contexts and scene complexities. For the experiments, 1,000 images were randomly selected from the validation set, ensuring a representative sample across different object categories, lighting conditions, scales, and occlusion scenarios.

For data processing, images were initially handled at their original resolution for object detection, then resized to 512x512 pixels during the editing phase to maintain compatibility with the Stable Diffusion model. The rich annotations of the dataset were utilized to validate the accuracy of both object detection and masking procedures. The selected images encompass a wide range of scenarios, including indoor and outdoor scenes, single and multiple object compositions, and varying levels of background complexity, providing a robust test bed for evaluating the model’s editing capabilities under real-world conditions.

3.2. Computing Resource

The experiments were conducted on a cloud computing platform using a high-performance server configuration. The system was equipped with an NVIDIA RTX A5000 GPU with 24GB GDDR6 memory, which provided sufficient computational power for the deep learning models including DETR, CLIP, and Stable Diffusion. The server featured an Intel Xeon Silver 4208 processor with 30 cores, supported by 118GB of RAM, ensuring efficient processing of large-scale image operations. Storage was configured with 145GB of disk space, and the system ran on Ubuntu 20.04 operating system. The network connection speed was maintained at 10000 Mb/s to facilitate smooth data transfer during the experiments. This hardware configuration enabled efficient processing of the image editing pipeline, with the GPU particularly crucial for the compute-intensive diffusion model operations.

3.3. Experimental design/setup

The experimental evaluation was designed to quantitatively measure the performance of the proposed image editing framework. Three key metrics were employed for performance assessment: Fréchet Inception Distance (FID)(Heusel et al., 2017) for measuring the quality and realism of edited images, Structural Similarity Index

Table 1. FID and SSIM score of each experiments

PROMPT	FID	SSIM	SUCCESS(%)
COLOR	36.23	0.6259 ± 0.1787	99.00
TILT	34.13	0.6323 ± 0.1719	98.90
FRONT	33.89	0.6306 ± 0.1760	99.00
UPDOWN	37.02	0.6307 ± 0.1797	99.40

(SSIM)(Wang et al., 2004) for evaluating the preservation of image structure outside the edited regions, and success rate representing the ratio of successfully completed edits to edit attempts. Additionally, processing time per image was recorded to assess practical applicability.

The experimental setup included four different types of object modifications: color transformation (e.g., "object changes to red object"), left-tilting (e.g., "object changes to object tilted slightly to the left"), front-facing orientation (e.g., "object changes to object is facing the front"), and upside-down orientation (e.g., "object changes to object is upside down"). These diverse modification tasks were designed to evaluate the model's capability in handling various types of object transformations while maintaining visual coherence. For each modification type, a total of 1,000 images were processed, with all metrics being recorded through an automated evaluation system. The testing process was designed to ensure reproducibility, with performance measurements taken at each stage of the pipeline from object detection through final editing. All experimental phases were conducted in the same computing environment to ensure consistent performance measurements.

3.4. Results

The proposed framework was evaluated across four different object modification tasks using a set of 1,000 images each from the COCO validation dataset. Each task demonstrated stable performance across all evaluation metrics with minor variations. Color transformation achieved an FID score of 36.23 and an average SSIM of $0.6259 (\pm 0.1787)$, with a notably high success rate of 99.00%. The left-tilting modification showed slightly improved results with an FID of 34.13 and SSIM of $0.6323 (\pm 0.1719)$, maintaining a high success rate of 98.90%. Front-facing orientation modification demonstrated the best FID score of 33.89 and a comparable SSIM of $0.6306 (\pm 0.1760)$, with a 99.00% success rate. The upside-down orientation task, while showing a slightly higher FID of 37.02, maintained consistent SSIM performance at $0.6307 (\pm 0.1797)$ and achieved the highest success rate at 99.40%. Processing times were remarkably consistent across all tasks, averaging between 2646.79ms and 2658.32ms per image, indicating stable computational performance regardless of the modification type. The consistency in SSIM scores (all approximately 0.63) suggests that

the framework maintains similar levels of structural preservation across different types of modifications. The low FID scores (ranging from 33.89 to 37.02) indicate high-quality, realistic results across all transformation types.

3.5. Discussion

The proposed framework demonstrates competitive performance when compared to existing image editing models. While previous approaches like DiffEdit(Couairon et al., 2022) achieved FID scores around 38-40 and SSIM values of approximately 0.60 for object-specific editing tasks, our method shows improved performance across all modification types, with FID scores ranging from 33.89 to 37.02 and SSIM values consistently above 0.63.

Particularly noteworthy is the model's performance in front-facing orientation modifications, which achieved the best FID score of 33.89, surpassing the previous state-of-the-art score of 35.5 reported by InstructPix2Pix(Brooks et al., 2022). The consistent SSIM scores (approximately 0.63) across all modification types also represent an improvement over existing methods, which typically show more variance between different types of edits. For instance, previous methods like SDEdit reported SSIM scores varying from 0.58 to 0.62 depending on the editing task.

However, the current approach has several limitations, particularly in the object masking process. For complex shapes or overlapping objects, accurate mask generation can be challenging. The bounding box-based masking generated by DETR fails to capture fine object boundaries, which can be problematic especially when editing objects with complex shapes or textures. Additionally, when multiple objects overlap, accurate mask separation becomes difficult, potentially leading to unintended regions being edited together. Nevertheless, the high success rates (98.90% - 99.40%) across all modification types represent a notable improvement compared to existing approaches, which typically achieve success rates between 95-97%. This improvement can be attributed to the effective integration of DETR's precise object detection with CLIP's semantic understanding. The processing time, while slightly longer than some existing methods, is offset by the improved quality and reliability of the edits.

4. Future Direction

This framework can be enhanced in two primary directions. First, the integration of instance segmentation models such as Mask2Former or SAM (Segment Anything Model) could replace the current bounding box-based masking approach, potentially leading to more precise object masks and better handling of complex shapes and overlapping objects. Second, the implementation of a more sophisticated CLIP score calculation method that considers multiple image patches

and their spatial relationships could improve the accuracy of object-text matching, particularly for complex scenes with multiple similar objects.

References

- Brooks, T., Holynski, A., and Efros, A. A. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23556–23565, 2022.
- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pp. 213–229. Springer, Cham, 2020.
- Couairon, G., Grechka, A., Verbeek, J., Schwenk, H., and Cord, M. Diffedit: Diffusion-based semantic image editing with mask guidance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9213–9223, 2022.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., and Zitnick, C. L. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pp. 740–755. Springer, Cham, 2014.
- Radford, A., Kim, J. H., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., and Sutskever, I. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- Wang, Z., Bovik, A. C., Sheikh, H. R., and Simoncelli, E. P. Image quality assessment: From error visibility to structural similarity. volume 13, pp. 600–612, 2004.