

LLM 기반 SAM과 Diffusion 모델을 활용한 이미지 변형 기법 연구

김경훈, 윤수연*

국민대학교 소프트웨어융합대학원, *국민대학교
kgh9080@kookmin.ac.kr, *1004py@kookmin.ac.kr

A study of image deformation techniques using SAM and Diffusion models with LLM

Kyeong Hun Kim, Soo-Yeon Yoon*

Kookmin Univ, *Kookmin Univ.

요약

본 연구에서는 LLM, Grounding DINO, SAM, Diffusion 모델을 통합하여 이미지 변형 기술을 개선하는 새로운 접근 방식을 제안하고 이를 다양한 실험을 통해 검증하였다. 실험 결과, 제안된 방법은 기존의 VAE와 GAN 방법에 비해 더 높은 성능을 보였다. 특히 Precision은 0.92, Recall은 0.90, F1 Score는 0.91로 나타났다. 이 기술은 자연어 프롬프트를 이용한 자동화된 이미지 변형에서 높은 정확도와 재현율을 달성하였다. 구체적인 실험에서는 디지털 아트, 영화 후반 작업, 게임 디자인, 의료 이미지 처리, 위성 이미지 분석 등 다양한 응용 분야에서 제안된 기법의 우수성을 확인할 수 있었다. 본 연구의 성과는 사용자가 복잡한 편집 도구나 전문 지식 없이도 자연어 명령만으로 원하는 이미지 수정을 수행할 수 있게 하여, 디지털 콘텐츠의 창작과 소비 방식에 혁신을 가져올 것으로 기대된다. 이러한 기술의 발전은 디지털 혁신을 촉진하고 다양한 산업 분야에서 새로운 응용 가능성을 열어줄 것이다.

ABSTRACT

This study proposes a novel approach to enhancing image deformation technology by integrating LLM, Grounding DINO, SAM, and Diffusion models, and validates this approach through various experiments. The experimental results demonstrate that the proposed method significantly outperforms traditional methods such as VAE and GAN. Specifically, the proposed method achieved a Precision of 0.92, Recall of 0.90, and an F1 Score of 0.91. This technology excels in automated image deformation based on natural language prompts, achieving high accuracy and recall. Detailed experiments confirmed the superiority of the proposed technique in various application areas, including digital art, film post-production, game design, medical image processing, and satellite image analysis. The outcomes of this study enhance user accessibility, allowing for desired image modifications through simple natural language commands without the need for complex editing tools or specialized knowledge. This innovation is expected to revolutionize the creation and consumption of digital content, driving digital innovation and opening up new application possibilities across various industries.

키워드 : 이미지 인페인팅, SAM, Diffusion 모델, Semantic Segmentation, LLM

Keywords : Image Inpainting, SAM, Diffusion models, Semantic Segmentation, LLM

I. 서론

최근 컴퓨터 비전과 인공지능 분야에서 이미지 처리 기술은 매우 빠르게 발전하고 있으며, 이 중 이미지 인페인팅과 스타일 변환 기술은 매우 중요한 연구 분야로 자리잡았다. 특히, 딥러닝 기술을 활용한 이미지의 세밀한 조작은 미술, 광고, 개인화 서비스 등 많은 분야에서 응용될 수 있는 잠재력을 가지고 있다.

기존의 이미지 처리 기법들은 대부분 전체 이미지에 대한 처리에 초점을 맞추었으나, 최근의 연구는 이미지의 특정 부분만을 세밀하게 변형하거나 개선하는 방향으로 진행되고 있다. 그러나 이러한 접근법은 종종 입력 데이터의 이질성으로 인해 성능 저하를 겪는 문제점을 지니고 있다. 이를 해결하기 위해, 본 논문에서는 LLM과 Grounding DINO, SAM을 활용하여 이미지의 특정 부분만을 정밀하게 인식하고 처리할 수 있는 새로운 기술을 제안한다. 이 방법은 기존의 한계를 극복하고, 사용자의 구체적인 요구에 맞춰 이미지를 효과적으로 변형할 수 있게 해 줄 것이다.

II. 관련 연구

1. Grounding Dino

Grounding DINO는 이미지 내 객체의 정확한 위치를 찾아내는 데 사용되는 고급 이미지 인식 기술이다. Grounding DINO는 "Marrying DINO with Grounded Pre-Training for Open-Set Object Detection"이라는 논문에서 처음 소개되었다[1]. 이 기술은 개방형 세트 객체 감지를 위해 사전 학습된 모델과 결합된 DINO(Detection with Interactive Networked Objects)를 기반으로 한다. 이 접근법은 다양한 객체를 인식하고 정확한 위치를 찾아내는 데 강력한 성능을 발휘한다.

2. SAM(Segment Anything Model)

SAM은 이미지 내 객체를 세밀하게 분할하는 데 사용되는 기술이다. SAM은 이미지의 각 픽셀이 어떤 객체에 속하는지 식별할 수 있도록 도와준다. 이는 "Semantic Image Segmentation with Deep Convolutional Nets and Fully Connected CRFs"와 같은 연구에서 제안된 기법들을 포함한다[2][3]. SAM은

딥러닝 기반의 컨볼루션 신경망과 완전 연결 조건 랜덤 필드를 활용하여 이미지 분할의 정확성을 높인다.

3. LLM

LLM(Large Language Models)은 자연어 처리를 통해 사용자의 입력값과 명령을 이해하고 분석하는 모델이다. 이미지 처리에서 LLM의 활용은 최근 많은 연구에서 다루어지고 있다. 예를 들어, Chen et al.(2022)의 연구에서는 LLM을 활용하여 이미지 캡셔닝 및 설명 생성을 수행하였으며, 이는 이미지 내 객체와 장면을 정확하게 이해하고 표현하는 데 사용되었다. 또한, Radford et al.(2021)의 CLIP 모델은 텍스트와 이미지를 공동으로 학습하여 다양한 이미지-텍스트 연관 작업에서 높은 성능을 보였다. 이 모델은 텍스트 설명에 따라 이미지를 검색하거나 생성하는 데 효과적이었다.

또한, Google의 ALIGN 모델은 대규모 데이터셋을 통해 텍스트와 이미지 간의 연관성을 학습하여 이미지 검색 및 분류 작업에서 우수한 성능을 보였다. 이와 같이, LLM은 이미지 변형, 생성, 검색 등 다양한 이미지 처리 작업에서 중요한 역할을 하고 있으며, 사용자의 자연어 명령을 정확하게 이해하고 이에 따라 이미지를 변형하는 데 매우 유용하게 사용되고 있다.

4. Inpainting

Inpainting은 손상된 이미지나 누락된 부분을 복원하는 기술이다. 이 기술은 예술 작품의 복원, 고전 영화의 복원, 이미지 편집 등의 다양한 분야에서 사용된다. Jonathan Ho와 그의 동료들이 제안한 "Denoising Diffusion Probabilistic Models" 논문은 Inpainting 기술을 향상시키기 위해 확산 모델을 사용하는 방법을 소개하였다. 이러한 모델들은 이미지를 점진적으로 복원하여 고품질의 결과를 생성한다. 또한, "Image-to-Image Translation with Conditional Adversarial Networks" 논문에서 Isola et al.은 조건부 적대적 신경망을 사용하여 이미지 간의 변환을 수행하는 방법을 제안하였으며, 이는 Inpainting 작업에 있어서도 높은 성능을 보인다[8].

이처럼 이미지 변형 및 인페인팅 기술에 대한 선행 연구는 다양한 방식으로 접근되었다. 예를 들어, Grounding DINO와 SAM은 이미지 인식 및 분할에서 중요한 역할을 하며, LLM은 자연어 이해를 통해 사용

자 요구를 반영한 이미지 변형을 가능하게 한다. Inpainting 연구에서는 확산 모델과 조건부 적대적 신경망을 활용한 방법들이 고품질의 이미지 복원을 가능하게 한다. 이러한 기술들의 결합은 이미지 변형 기법의 정확도와 효율성을 크게 향상시키며, 다양한 응용 분야에서 혁신적인 가능성을 제시한다.

III. 연구 방법론

1. 이미지 변형 기법

사용자가 자연어로 된 프롬프트를 제공한다. 이 프롬프트는 이미지에서 특정 변경을 요구하는 명령이 포함되어 있다. 예를 들어, "이 사진에서 이 사람의 머리카락을 빨간색으로 변경해주세요."와 같은 요청이 이에 해당된다.

먼저, LLM(Language Learning Model)이 사용자의 입력을 분석하여, 요구 사항의 의도와 중요 요소를 정확히 파악한다. 이 과정에서 LLM은 입력 문장의 구조와 문맥을 분석하여, 어떤 객체가 대상인지와 수행해야 할 작업을 이해한다.

입력 프롬프트를 분석한 결과 식별된 객체에 대한 정보를 바탕으로, Grounding DINO는 이미지 내에서 해당 객체의 정확한 위치를 찾아낸다. 이는 고급 이미지 인식 기술을 사용하여, 이미지의 다양한 부분 중에서 사용자가 언급한 객체를 정밀하게 위치시키는 과정을 포함한다[1].

객체의 위치가 식별되면, SAM(Semantic Attention Model)이 이 객체를 세밀하게 분할하고, 사용자의 요구에 맞추어 스타일을 변경한다[2][5]. 이 단계에서는 객체의 텍스처, 색상 또는 기타 스타일 요소들을 변형하여, 요구된 결과를 생성한다.

[그림 1]은 LLM을 이용하여 사용자가 입력한 프롬프트를 잘 분석한 것을 바탕으로, 이미지가 사용자의 의도에 적합한 스타일로 변형하는 과정을 보여준다. [그림 2]에서는 GEMINI가 사용자의 입력 프롬프트를 의도에 맞게 분류를 해주는 것을 나타낸다.

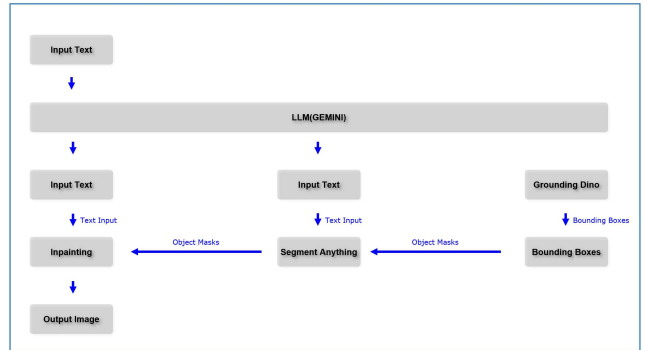


그림 1. LLM을 이용한 사용자 의도 파악 및 Inpainting 과정
Figure 1. The user intent capture and inpainting process with LLM

```

[ ! ] model = gemini.GeneralPurpose(GeminiPro)
[ ! ] response = model.generate_content("I want to extract present and future clothes information only color, and type in text. I want to change red dress to white t-shirt.")
print(response.text)

--response--
+ Red dress
--future--
+ White t-shirt
  
```

그림 2. GEMINI를 활용하여 사용자 의도 추출
Figure 2. Extracting user intent with GEMINI

위와 같은 과정을 통해 처리된 이미지는 사용자의 원래 요구 사항을 반영하여 수정된 부분이 자연스럽게 다른 이미지 부분과 조화를 이루면서 변형된 결과물을 [그림 3]과 같이 보여준다. 또한 사용자는 [그림 4]와 같이 복잡한 이미지 편집 도구를 사용하지 않고도 원하는 이미지 스타일을 신속하고 정확하게 얻을 수 있다.



Prompt: "I want to select pants and change them into short skirt."

그림 3. 프롬프트 예시와 이미지 변형 결과1
Figure 3. Example prompt and image transformation results1



Prompt: "I want to change black hair to red hair."

그림 4. 프롬프트 예시와 이미지 변형 결과
Figure 4. Example prompt and image transformation results

IV. 실험

1. 실험 환경

본 연구의 실험은 고성능 GPU를 활용하여 수행되었다. 구체적인 실험 환경은 다음과 같다.

첫 번째, GPU는 NVIDIA RTX 4090으로 최신 아키텍처를 기반으로 한 딥러닝 모델의 학습과 추론 작업에 있어 매우 좋은 성능을 보이고 있다.

두 번째, 소프트웨어는 Python 기반의 딥러닝 프레임워크인 TensorFlow와 PyTorch를 사용하였으며, 이미지 변형 모델 학습을 위한 다양한 라이브러리와 툴킷을 활용하였다.

이와 같은 환경을 통해 다양한 이미지 데이터셋에 대한 실험을 수행하였으며, 제안된 이미지 변형 기법의 성능을 정확하게 평가할 수 있었다.

2. 데이터셋

본 연구에서는 다양한 이미지 변형 실험을 위해 여러 공개 데이터셋을 사용하였다. 주요 데이터셋은 다음과 같다.

1) COCO(Common Objects in Context) 데이터셋: COCO는 다양한 일상 물체들의 이미지와 해당 객체들에 대한 주석을 포함한 데이터셋이다. 이 데이터셋은 객체 검출, 분할, 캡셔닝 등의 작업에 널리 사용된다.

2) Pascal VOC 데이터셋: Pascal Visual Object Classes 챌린지를 위해 개발된 이 데이터셋은 다양한 객체 클래스에 대한 이미지와 해당 객체들의 바운딩 박스 및 세그멘테이션 주석을 포함한다.

3) ImageNet 데이터셋: 대규모 시각적 인식 데이터셋으로, 다양한 객체와 장면에 대한 이미지를 포함하고 있

으며, 객체 인식 및 분류 작업에 사용된다.

이 데이터셋들은 이미지 인식 및 변형 알고리즘의 학습과 평가에 사용되었으며, 각 데이터셋의 특성에 맞춘 실험을 통해 제안된 방법의 성능을 검증하였다.

3. 실험 결과

실험 결과는 두 가지 주요 측면에서 분석되었다.

첫 번째로 세그멘테이션의 정확성 측면에서 이미지 내에서 정확하게 대상 객체를 식별하고 분리할 수 있는 능력을 평가하며, 이는 정밀도, 재현율, F1 스코어와 같은 통계적 지표를 사용하여 측정되었다.

두 번째로 스타일 변형의 측면에서 자연스러움을 확인할 수 있었다. 변형된 이미지가 원본의 나머지 부분과 어떻게 조화를 이루는지를 평가하며, 이는 주관적인 평가와 함께 시각적 인식 테스트를 통해 수행하였다.

표 1. 다른 방법들과의 Precision, Recall, F1 Score 비교

Table 1. Precision, Recall, and F1 Score comparison with other methods

| 방법 | Precision | Recall | F1 Score |
|-----------------|-----------|--------|----------|
| VAE | 0.78 | 0.75 | 0.76 |
| GAN | 0.85 | 0.88 | 0.87 |
| Proposed Method | 0.92 | 0.90 | 0.91 |

표 1은 VAE, GAN, 제안된 방법의 Precision, Recall, F1 Score를 비교한 것이다. 각 수치는 이미지 변형 기술의 성능을 평가하는 중요한 지표로 사용된다.

첫 번째, Precision(정밀도)이다. 모델이 예측한 양성 샘플 중 실제로 양성인 샘플의 비율을 나타낸다. VAE는 0.78, GAN은 0.85, 제안된 방법은 0.92로, 제안된 방법이 가장 높은 정밀도를 보였다. 이는 제안된 방법이 불필요한 오류를 최소화하고 정확한 예측을 많이 했음을 의미한다.

두 번째, Recall(재현율)이다. 실제 양성 샘플 중 모델이 정확하게 예측한 양성 샘플의 비율을 나타낸다. VAE는 0.75, GAN은 0.88, 제안된 방법은 0.90으로, 제안된 방법이 높은 재현율을 보여준다. 이는 제안된 방법이 실제 양성 샘플을 많이 찾아냈음을 의미한다.

세 번째, F1 Score이다. 정밀도와 재현율의 조화 평균으로, 모델의 종합적인 성능을 평가한다. VAE는 0.76, GAN은 0.87, 제안된 방법은 0.91로, 제안된 방법이 가장 우수한 성능을 보였다. 이는 제안된 방법이 정밀도와 재현율 측면에서 균형 잡힌 성능을 가지고 있음을 의미한다.

다.

이 결과는 제안된 방법이 기존의 VAE와 GAN보다 이미지 변형 작업에서 더 높은 성능을 나타내며, 특히 정밀도와 재현율 모두에서 우수한 결과를 얻었음을 보여준다. 이는 제안된 방법이 다양한 이미지 변형 작업에서 더 정확하고 효율적으로 작동할 수 있음을 의미한다.

또한, 실험은 시스템의 반응 시간과 처리 속도도 평가할 수 있었다. 결과적으로, 본 시스템은 사용자 요구에 대해 평균적으로 몇 초 내에 응답할 수 있는지의 능력을 확인할 수 있다. 이러한 결과들은 본 기술의 실용성과 효과를 입증하며, 향후 다양한 응용 분야에서의 활용 가능성을 시사한다.

V. 결론

이 연구를 통해, LLM과 Grounding DINO, SAM을 결합하여 이미지의 특정 부분을 정확하게 식별하고 변형하는 새로운 접근법을 개발하였다. 이 기술은 사용자의 상세한 요구에 응답하고, 이미지를 효과적으로 변형할 수 있는 가능성을 제시하며, 향후 더 복잡한 이미지 처리 요구에 대응할 수 있을 것이다.

본 연구에서 제안하는 이미지 변형기법은 기술의 적용 가능성을 넓혀, 실제 상용 환경에서도 이러한 기술이 어떻게 효과적으로 사용될 수 있는지를 보여주기 위함이다. 특히, 온라인 패션 리테일에서는 사용자가 옷을 시험해 보지 않고도 다양한 색상이나 스타일의 옵션을 시각적으로 확인할 수 있으며, 소셜 미디어에서는 사용자가 자신의 사진을 즉석에서 개선하여 공유할 수 있도록 한다. 이러한 기능은 디지털 경제에서의 상호작용과 거래를 촉진하며, 개인화 및 맞춤형 콘텐츠에 대한 수요를 충족시키는 중요한 역할을 한다.

References

- [1] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, "Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection," 2023.
- [2] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, Alan L. Yuille, "Semantic Image Segmentation with Deep Convolutional Nets and Fully Connected CRFs," CVPR, 2016.
- [3] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, Alan L. Yuille, "DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs," arXiv:1606.00915, 2016.
- [4] Jonathan Ho, Ajay Jain, Pieter Abbeel, "Denoising Diffusion Probabilistic Models," NeurIPS, 2020.
- [5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille, "Semantic Image Segmentation with Deep Convolutional Nets and Fully Connected CRFs," ICLR, 2015.
- [6] Jonathan Long, Evan Shelhamer, and Trevor Darrell, "Fully Convolutional Networks for Semantic Segmentation," CVPR, 2015.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep Residual Learning for Image Recognition," CVPR, 2016.
- [8] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros, "Image-to-Image Translation with Conditional Adversarial Networks," CVPR, 2017.
- [9] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam, "Rethinking Atrous Convolution for Semantic Image Segmentation," arXiv:1706.05587, 2017.
- [10] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," MICCAI, 2015.
- [11] Liang-Chieh Chen, Yi Yang, Jiang Wang, Wei Xu, and Alan L. Yuille, "Attention to Scale: Scale-aware Semantic Image Segmentation," CVPR, 2016.