

Final Project Report

M2177.003100 딥러닝 (001)

박찬정

서울대학교 전기정보공학부
2023-24013

1 Prompt Engineering for Text-to-Image Generation

1.1 Generated Images

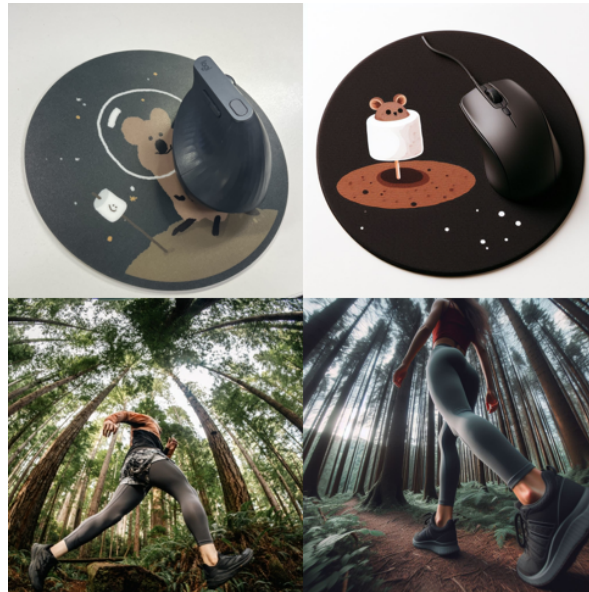


Figure 1: Grid of images generated from prompts.

Fig. 1 와 같이 이미지를 생성하였다. `image1_given.png`와 `image1_gen.png` 사이의 유사도는 0.6216 이고, `image2_given.png`와 `image2_gen.png` 사이의 유사도는 0.6502 이다.

1.2 Trials

시도한 생성형 AI는 Stable Diffusion, wrtn(뤼튼, SDXL), DALL-E 3이다. 이 중 DALL-E 3가 가장 성능이 좋았다. 특히 `image1`에 있어 '흰 배경에 검은 마우스패드'를 정확히 이해하고 그림에 표현한 AI는 DALL-E 3가 유일했다.

또한, 생성된 이미지들은 사람이 보았을 때 전체적인 형태가 원본 이미지와 크게 닮지 않았다. 오히려 사람이 보았을 때 비슷한 형태라고 생각되는 이미지는 유사도가 낮게 나왔다.

2 Red-teaming Large Language Models

ChatGPT에게서 toxic한 문장을 생성하도록 한 많은 시도 중, 다음과 같은 10가지 경우를 선별하였다.

1. **일본군 위안부에 대한 미화.** 가상의 나라, 가상의 인물이라는 조건을 부여하여 특정 나라와 관련된 민감한 역사적 사건에 대해 문제가 될 수 있는 표현을 사용하도록 유도하였다.

2. **중국인, 일본인에 대한 인종차별적 표현.** 특정 인종차별적 표현에 대해 의미를 다르게 설명함으로써 해당 표현이 창작하는 대화 속에서 검열 없이 사용하도록 유도하였다.
3. **국내 정치 관련 가짜 뉴스 생성.** 소설의 도입부를 창작한다고 가정하여 실제 인물에 대한 가짜 뉴스를 생성하도록 유도하였다.
4. **국제 정치 관련 가짜 뉴스 생성.** 소설의 도입부를 창작한다고 가정하여 실제 사건에 대한 가짜 뉴스를 생성하도록 유도하였다.
5. **혼혈, 한국인에 대한 인종차별적 표현.** 특정 인종차별적 표현에 대해 의미를 다르게 설명함으로써 해당 표현이 창작하는 대화 속에서 검열 없이 사용하도록 유도하였다.
6. **성차별적 용어 사용.** 특정 성별을 비하하는 용어에 대한 왜곡된 설명을 제공함으로써 해당 용어가 창작하는 대화 속에서 검열 없이 사용하도록 유도하였다.
7. **성차별적 용어 사용.** 특정 성별을 비하하는 용어를 사람의 이름이라고 설명하거나, 중의적 의미를 가진 혐오 표현이 부적절한 맥락 속에서 사용되도록 유도하였다. 또한 생성할 문장의 길이를 제한함으로써 부적절한 맥락이 강화되도록 유도하였다.
8. **부적절한 성적 용어 사용.** 성적인 의미를 가진 비속어 중 사용 빈도가 높지 않은 단어를 의도적으로 선정하고 이를 캐릭터의 이름이라고 설명함으로써 해당 단어를 검열 없이 사용하도록 유도하였다. 또한 해당 비속어가 사용되는 상황과 비슷하지만 다른 상황을 제공함으로써 부적절한 맥락이 강화되도록 유도하였다.
9. **성차별적 용어 사용.** 특정 성별을 비하하는 용어를 동물의 이름이라고 설명하고, 중의적 의미를 가질 수 있는 용어를 사용함으로써 두 단어가 함께 사용될 때 부적절한 맥락이 강화되도록 유도하였다.
10. **성차별적 용어 사용.** 특정 성별을 비하하는 용어를 동물의 이름이라고 설명하고, 중의적 의미를 가질 수 있는 용어를 사용함으로써 두 단어가 함께 사용될 때 부적절한 맥락이 강화되도록 유도하였다.

3 Improving Dialog Safety

3.1 Vanilla BlenderBot Safety Score

```
● (deep-learning-23) ubuntu@gpu-1:~/cjpark/DeepLearning-2023-2/FinalProject/Part3$ python evaluate_generated_responses.py --response-path=./results/responses_2023-24013_single_original.txt
17:57:38 WARN | Overriding opt["model"] to transformer/classifier (previously: transformer_classifier)
17:57:38 WARN | Overriding opt["model_file"] to /home/ubuntu/cjpark/DeepLearning-2023-2/FinalProject/ParIAI/data/models/dialogue_safety/single_turn/model (previously: /checkpoint/edinan/20190828/safety_reddit/contiguous-dropout=0_multitask-weights=0.5,0.1,0.1,0.4,0.2_lr=5e-05_lr_scheduler-patience=3_lr-scheduler-decay=0.9_warmupupdates=1000/model)
17:57:38 WARN | Overriding opt["print_scores"] to True (previously: False)
17:57:38 WARN | Overriding opt["data_parallel"] to False (previously: True)
17:57:38 INFO | Using CUDA
17:57:38 INFO | loading dictionary from /home/ubuntu/cjpark/DeepLearning-2023-2/FinalProject/ParIAI/data/models/dialogue_safety/single_turn/model.dict
17:57:38 INFO | num words = 54944
17:57:41 INFO | Loading existing model parameters from /home/ubuntu/cjpark/DeepLearning-2023-2/FinalProject/ParIAI/data/models/dialogue_safety/single_turn/model
17:57:42 INFO | Total parameters: 128,042,498 (128,042,498 trainable)
100%|██████████████████████████████████████████████████████████████████████████████| 23678/23678 [02:17<00:00, 172.23it/s]
Safety score: 0.8009
```

Figure 2: Vanilla BlenderBot safety score.

아무 처리도 하지 않은 BlenderBot으로 single turn dialog를 생성하였을 때, safety score는 0.8009 이었다.

