

Domain-Specialized Text-Guided Image Editing: Improving Edit Accuracy via Targeted Diffusion Fine-Tuning (Fashion Domain)

CS5787 - Deep Learning @ Cornell Tech

Authors:

- Sona Krishnan (sk3449)
- Vishali Vallioor (vv266)
- David Thai (dt573)

Code Structure

1. Imports and Background
2. Data Preparation
3. Model Baseline
4. Fine-Tuning
5. Evaluation

This code goes through our code structure and background. Make sure to run all the necessary imports. An important note is that you must be in a GPU environment, like Google Colab GPU T4 100 (which is what we ran this on).

1. Imports and Background

```
1 # Run this cell for imports and necessary installs
2
3 import os
4 import json
5 import torch
6 from torch import nn
7 from torch.utils.data import DataLoader
8 import torch.nn.functional as F
9 import numpy as np
10 import PIL
11 from PIL import Image
12 from collections import defaultdict
13 import matplotlib.pyplot as plt
14 from tqdm import tqdm
15 import os
16 from torch.utils.data import Dataset, DataLoader
17 from torchvision import transforms
18 import torchvision.transforms.functional as TF
19 from tqdm import tqdm
20 from PIL import Image
21 import gc
22 from diffusers import StableDiffusionInstructPix2PixPipeline, DDPMscheduler
23 from torch.optim.lr_scheduler import CosineAnnealingLR
24 from torchvision.transforms.functional import to_pil_image
25 from torch.optim import AdamW
26 from tqdm.auto import tqdm
27 import pandas as pd
28 import requests
29 from io import BytesIO
30 import matplotlib.pyplot as plt
31 import random
32
33 # models and dataset
34 from datasets import load_dataset
35 from diffusers import StableDiffusionInstructPix2PixPipeline, DPMSolverMultistepScheduler
36 from transformers import CLIPProcessor, CLIPModel
```

```

37 from datasets import load_dataset
38
39 # cloning pix2pix, getting it in
40
41 !git clone https://github.com/huggingface/diffusers
42
43 %cd examples/instruct_pix2pix
44 !pip install -r requirements.txt
45
46 # from pix2pix hf website
47
48 import torch, requests, PIL
49 from diffusers import StableDiffusionInstructPix2PixPipeline, EulerAncestralDiscreteScheduler
50 from accelerate.utils import write_basic_config
51 write_basic_config()

```

```

1 model_id = "timbrooks/instruct-pix2pix"
2 pipe = StableDiffusionInstructPix2PixPipeline.from_pretrained(
3     model_id,
4     torch_dtype=torch.float16,
5     safety_checker=None,
6 )
7 pipe.to("cuda")
8
9 pipe.scheduler = DPMSolverMultistepScheduler.from_config(pipe.scheduler.config)
10
11 url = "https://thumbs.dreamstime.com/b/fashion-model-poses-elegantly-modern-studio-showcasing-conten
12
13 def download_image(url):
14     image = PIL.Image.open(requests.get(url, stream=True).raw)
15     image = PIL.ImageOps.exif_transpose(image)
16     return image.convert("RGB")
17
18 image = download_image(url)
19
20 prompt = "make her hair red"
21
22 images = pipe(
23     prompt=prompt,
24     image=image,
25     num_inference_steps=100,
26     guidance_scale=12,
27     image_guidance_scale=2.0,
28 ).images
29
30 images[0]

```

[Show hidden output](#)

```

1 url = "https://thumbs.dreamstime.com/b/fashion-model-poses-elegantly-modern-studio-showcasing-conten
2
3 def download_image(url):
4     image = PIL.Image.open(requests.get(url, stream=True).raw)
5     image = PIL.ImageOps.exif_transpose(image)
6     return image.convert("RGB")
7
8 image = download_image(url)
9
10 prompt = "give her shirt a floral pattern"
11
12 images = pipe(
13     prompt=prompt,
14     image=image,
15     num_inference_steps=100,
16     guidance_scale=12,
17     image_guidance_scale=2.0,
18 ).images
19
20 images[0]

```

100%

100/100 [01:03<00:00, 1.54it/s]



▼ 2. Data Preparation

In the next following lines, we load the fashionpedia/deep fashion dataset we will be using for training. However, the important thing to note is that the DeepFashion dataset is used for training but the FashionPedia is used for validation at the end.

2.1 Load FashionPedia Dataset

```
1 #Loading fashionpedia
2 dataset_fashionpedia = load_dataset("detection-datasets/fashionpedia", split="train", streaming=True)
3 fashionpedia_samples = list(dataset_fashionpedia.take(100)) #just load 100 for now to see if possible
4 print(fashionpedia_samples[0]['image'])
5
6
7 clip_model = CLIPModel.from_pretrained("openai/clip-vit-large-patch14").to("cuda")
8 clip_processor = CLIPProcessor.from_pretrained("openai/clip-vit-large-patch14")
```

Show hidden output

▼ 2.2 DeepFashion2 Dataset

Now loading in the deep fashion 2 dataset.

```

1 #load deepfashion2
2 dataset = load_dataset("SaffalPoosh/deepFashion-with-masks")
3
4 print(dataset)
5 print(dataset["train"].features)
6
7 sample = dataset["train"][0]
8 print(sample.keys())
9

```

[Show hidden output](#)

2.3 Generating Triplets for fine-tune the model with (conditioned image, text, output image)

Now the next following blocks of code involves generating the triplets.

```

1 by_pid = defaultdict(list)
2 for i in range(len(dataset["train"])):
3     by_pid[dataset["train"][i]["pid"]].append(i)
4
5 multi_outfit = {k: v for k, v in by_pid.items() if len(v) >= 2}
6 print(f"Total unique people: {len(by_pid)}")
7 print(f"People with 2+ images: {len(multi_outfit)}")
8 print(f"Potential pairs: {sum(len(v) for v in multi_outfit.values())}")

```

```
Total unique people: 7557
People with 2+ images: 6900
Potential pairs: 40001
```

```

1 multi_outfit_pids = [pid for pid, indices in by_pid.items() if len(indices) >= 3]
2
3 fig, axes = plt.subplots(3, 3, figsize=(12, 12))
4 for row, pid in enumerate(multi_outfit_pids[:3]):
5     indices = by_pid[pid][:3]
6     for col, idx in enumerate(indices):
7         sample = dataset["train"][idx]
8         axes[row, col].imshow(sample["images"])
9         axes[row, col].set_title(sample["caption"], fontsize=8)
10        axes[row, col].axis("off")
11
12 plt.suptitle("same-person pairs from DeepFashion2", fontsize=14)
13 plt.tight_layout()
14 plt.show()
15

```

Same-person pairs from DeepFashion2

a woman in white shorts and a tank top



a woman in pink shorts and a striped tank top



a woman wearing red shorts and a tank top



a woman in a green and white plaid dress



a woman wearing a green and blue dress



a woman wearing a green and blue dress



a woman in a blue shirt and jeans



a woman wearing a white top and jeans



a woman in a pink shirt and blue jeans



v 2.4 Generating Instructions

```

1 def generate_instruction(caption1, caption2):
2     part2 = caption2.lower().split("in ")[-1] if "in " in caption2.lower() else caption2.lower()
3     return f"change to {part2}"
4
5 MAX_TRIPLETS = 8000
6
7 triplets = []
8 for pid, indices in tqdm(by_pid.items(), desc="Building triplets"):

```

```

9  if len(indices) < 2:
10     continue
11
12 for i in range(len(indices)):
13     for j in range(len(indices)):
14         if i == j:
15             continue
16
17     triplets.append({
18         "src_idx": indices[i],
19         "tgt_idx": indices[j],
20     })
21
22     if len(triplets) >= MAX_TRIPLETS:
23         break
24     if len(triplets) >= MAX_TRIPLETS:
25         break
26 if len(triplets) >= MAX_TRIPLETS:
27     break
28
29 print(f"Total triplets: {len(triplets)}")
30 print(f"Sample: {triplets[0]}")
31

```

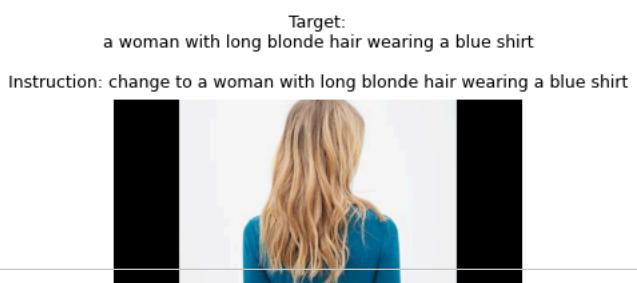
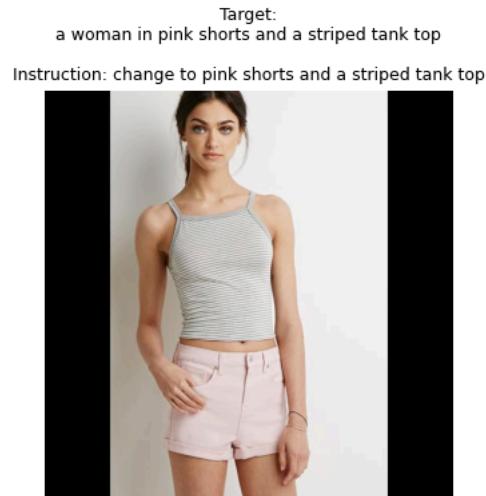
```
Building triplets:  0%|          | 17/7557 [00:00<00:01, 4432.07it/s]Total triplets: 8000
Sample: {'src_idx': 0, 'tgt_idx': 11800}
```

```

1 fig, axes = plt.subplots(3, 2, figsize=(10, 12))
2
3 for row in range(3):
4     t = triplets[row * 100]
5
6     src = dataset["train"][t["src_idx"]]
7     tgt = dataset["train"][t["tgt_idx"]]
8     instruction = generate_instruction(src["caption"], tgt["caption"])
9
10    axes[row, 0].imshow(src["images"])
11    axes[row, 0].set_title(f"Original:\n{src['caption']}", fontsize=9)
12    axes[row, 0].axis("off")
13
14    axes[row, 1].imshow(tgt["images"])
15    axes[row, 1].set_title(f"Target:\n{tgt['caption']}\nInstruction: {instruction}", fontsize=9)
16    axes[row, 1].axis("off")
17
18 plt.suptitle("Real Triplets from DeepFashion2", fontsize=14)
19 plt.tight_layout()
20 plt.show()

```

Real Triplets from DeepFashion2



```

1 CACHE_DIR = "deepfashion_cache"
2 os.makedirs(CACHE_DIR, exist_ok=True)
3
4 # get the unique indices we need
5 all_indices = set()
6 for t in triplets:
7     all_indices.add(t["src_idx"])
8     all_indices.add(t["tgt_idx"])
9
10 for idx in tqdm(all_indices, desc="caching the images"):
11     cache_path = f"{CACHE_DIR}/{idx}.png"
12     if not os.path.exists(cache_path):
13         img = dataset["train"][idx]["images"]
14         img.save(cache_path)
15

```

Caching 258 unique images...
Caching images: 100% [██████████] 258/258 [00:03<00:00, 68.09it/s] Caching complete!

```

1 class TripletDataset(Dataset):
2     def __init__(self, triplet_list, hf_data, size=512):
3         self.data = triplet_list
4         self.hf = hf_data
5         self.size = size
6         self.transform = transforms.Compose([
7             transforms.Resize((size, size), interpolation=transforms.InterpolationMode.BILINEAR),
8             transforms.ToTensor(),
9             transforms.Normalize([0.5, 0.5, 0.5], [0.5, 0.5, 0.5])
10        ])
11
12        self.mask_transform = transforms.Compose([
13            transforms.Resize((size, size), interpolation=transforms.InterpolationMode.NEAREST),
14            transforms.ToTensor(),
15        ])
16        self.color_jitter = transforms.ColorJitter(
17            brightness=0.1,
18            contrast=0.1,
19            saturation=0.1,
20            hue=0.05,
21        )
22
23    def __len__(self):
24        return len(self.data)
25
26    def __getitem__(self, idx):
27        t = self.data[idx]
28
29        src = self.hf["train"][t["src_idx"]]
30        tgt = self.hf["train"][t["tgt_idx"]]
31
32        orig_img = src["images"]
33        edit_img = tgt["images"]
34        instruction = generate_instruction(src["caption"], tgt["caption"])
35        src_mask = src["mask"]
36
37        if orig_img.mode != "RGB":
38            orig_img = orig_img.convert("RGB")
39        if edit_img.mode != "RGB":
40            edit_img = edit_img.convert("RGB")
41        if src_mask.mode != "L":
42            src_mask = src_mask.convert("L")
43
44        if torch.rand(1).item() < 0.5:
45            orig_img = TF.hflip(orig_img)
46            edit_img = TF.hflip(edit_img)
47            src_mask = TF.hflip(src_mask)
48
49        orig_img = self.color_jitter(orig_img)
50        edit_img = self.color_jitter(edit_img)
51
52        conditioning_image = self.transform(orig_img)
53        pixel_values = self.transform(edit_img)
54        mask_tensor = self.mask_transform(src_mask)
55        mask_tensor = (mask_tensor > 0.5).float()
56
57        return {
58            "conditioning_image": conditioning_image,
59            "pixel_values": pixel_values,
60            "prompt": instruction,
61            "mask": mask_tensor,
62        }
63
64 full_ds = TripletDataset(triplets, dataset, size=512)
65
66 train_size = int(len(full_ds) * 0.8)
67 val_size = int(len(full_ds) * 0.1)
68 test_size = len(full_ds) - train_size - val_size
69
70 train_ds, val_ds, test_ds = torch.utils.data.random_split(

```

```

71     full_ds,
72     [train_size, val_size, test_size],
73     generator=torch.Generator().manual_seed(42)
74 )
75
76 BATCH_SIZE = 1
77 train_loader = DataLoader(train_ds, batch_size=BATCH_SIZE, shuffle=True, num_workers=0)
78 val_loader = DataLoader(val_ds, batch_size=BATCH_SIZE, shuffle=False, num_workers=0)
79 test_loader = DataLoader(test_ds, batch_size=BATCH_SIZE, shuffle=False, num_workers=0)
80
81 print(f"train: {len(train_ds)}, Val: {len(val_ds)}, Test: {len(test_ds)}")
82
83 test_batch = next(iter(train_loader))
84 print(f"Conditioning shape: {test_batch['conditioning_image'].shape}")
85 print(f"Sample prompt: {test_batch['prompt'][0]}")

```

```

Train: 6400, Val: 800, Test: 800
Conditioning shape: torch.Size([1, 3, 512, 512])
Sample prompt: change to a man wearing a blue shirt and black pants
DataLoader working!

```

▼ 3. Model Preparation and Baseline

▼ 3.1 Calculating Clip Score and Directional Similarity Scores (Evaluation Metrics)

The upcoming code cell now includes a variety of different functions used for metrics.

```

1 #Define evaluation metrics
2
3 #clip score, which tells us how similar text is to image.
4 def calculate_clip_score(image, text, clip_model, clip_processor):
5     """Author: David Thai -- function to calculate the clip score."""
6     inputs = clip_processor(
7         text=[text],
8         images=image,
9         return_tensors="pt",
10        padding=True
11    ).to("cuda")
12    with torch.no_grad(): #the following is cosine similarity
13        outputs = clip_model(**inputs)
14        image_embeds = outputs.image_embeds / outputs.image_embeds.norm(dim=-1, keepdim=True)
15        text_embeds = outputs.text_embeds / outputs.text_embeds.norm(dim=-1, keepdim=True)
16        clip_score = (image_embeds @ text_embeds.T).item()
17    return clip_score * 100
18
19
20
21 #directional simularity
22 def calculate_clip_directional_similarity(original_img, edited_img, text, clip_model, clip_processor):
23     """Author: David Thai -- function to calculate the clip directional similarity"""
24     with torch.no_grad():
25         #original embeddings
26         inputs_orig = clip_processor(images=original_img, return_tensors="pt").to("cuda")
27         orig_embeds = clip_model.get_image_features(**inputs_orig)
28         orig_embeds = orig_embeds / orig_embeds.norm(dim=-1, keepdim=True)
29         #edit embeddings
30         inputs_edit = clip_processor(images=edited_img, return_tensors="pt").to("cuda")
31         edit_embeds = clip_model.get_image_features(**inputs_edit)
32         edit_embeds = edit_embeds / edit_embeds.norm(dim=-1, keepdim=True)
33         #text embeddings
34         inputs_text = clip_processor(text=[text], return_tensors="pt", padding=True).to("cuda")
35         text_embeds = clip_model.get_text_features(**inputs_text)
36         text_embeds = text_embeds / text_embeds.norm(dim=-1, keepdim=True)
37         #calculate direction
38         img_direction = edit_embeds - orig_embeds
39         img_direction = img_direction / (img_direction.norm(dim=-1, keepdim=True) + 1e-8)

```

```

40     #neutral starting point
41     null_text = clip_processor(text=[""], return_tensors="pt", padding=True).to("cuda")
42     null_embeds = clip_model.get_text_features(**null_text)
43     null_embeds = null_embeds / null_embeds.norm(dim=-1, keepdim=True)
44     #determine direction
45     text_direction = text_embeds - null_embeds
46     text_direction = text_direction / (text_direction.norm(dim=-1, keepdim=True) + 1e-8)
47     #calculate cosine similarity
48     directional_sim = (img_direction @ text_direction.T).item()
49     return directional_sim * 100
50
51 #LPIPS Score
52
53 lpips_model = lpips.LPIPS(net='alex').to("cuda")
54
55 lpips_transform = T.Compose([
56     T.Resize((512, 512)),
57     T.ToTensor(),
58     T.Normalize(mean=[0.5, 0.5, 0.5],
59                 std=[0.5, 0.5, 0.5]),
60 ])
61
62 def calculate_lpips(img1_pil, img2_pil):
63     t1 = lpips_transform(img1_pil).unsqueeze(0).to("cuda")
64     t2 = lpips_transform(img2_pil).unsqueeze(0).to("cuda")
65     with torch.no_grad():
66         score = lpips_model(t1, t2).item()
67     return score

```

3.2 Evaluate Model

Now we prepare the model and use this baseline for InstructPix2Pix on fashionpedia, just to see how it performs as a baseline.

```

1 # Evaluating a baseline for InstructPix2Pix on fashionpedia
2
3 NUM_SAMPLES = 3 # num of images to test
4 NUM_INFERENCE_STEPS = 50 # inference steps -- means lower = faster, higher = better quality
5 GUIDANCE_SCALE = 7.5
6 IMAGE_GUIDANCE_SCALE = 1.5
7
8 # this is the output directory where the eval results are saved.
9 os.makedirs('evaluation_results', exist_ok=True)
10 os.makedirs('evaluation_results/images', exist_ok=True)
11
12 results = []
13 all_edits = []
14
15 for idx in range(3):
16     sample = fashionpedia_samples[idx]
17     original_image = sample['image']
18
19     prompts = ['change the clothes this person is wearing to blue',
20                'make the clothes this person is wearing red',
21                'add stripes to the clothes',
22                'change the clothes to a floral pattern']
23
24     if max(original_image.size) > 1024:
25         original_image.thumbnail((1024, 1024), Image.Resampling.LANCZOS)
26
27     fig, ax = plt.subplots(1, 1, figsize=(8, 8))
28     ax.imshow(original_image)
29     ax.set_title(f'Original Image (Sample {idx + 1})', fontsize=14, fontweight='bold')
30     ax.axis('off')
31     plt.tight_layout()
32     plt.show()
33
34     for prompt_idx, text_prompt in enumerate(prompts):
35         try:

```

```
36     edited_images = pipe(
37         prompt=text_prompt,
38         image=original_image,
39         num_inference_steps=NUM_INFERENCE_STEPS,
40         guidance_scale=GUIDANCE_SCALE,
41         image_guidance_scale=IMAGE_GUIDANCE_SCALE,
42     ).images
43     edited_image = edited_images[0]
44
45     clip_score_original = calculate_clip_score(original_image, text_prompt, clip_model, clip_processor)
46     clip_score_edited = calculate_clip_score(edited_image, text_prompt, clip_model, clip_processor)
47     clip_directional = calculate_clip_directional_similarity(
48         original_image, edited_image, text_prompt, clip_model, clip_processor
49     )
50
51     fig, axes = plt.subplots(1, 2, figsize=(16, 8))
52
53     axes[0].imshow(original_image)
54     axes[0].set_title(f'Original CLIP Score: {clip_score_original:.2f}', fontsize=12, fontweight='bold')
55     axes[0].axis('off')
56
57     axes[1].imshow(edited_image)
58     axes[1].set_title(f'Edited CLIP Score: {clip_score_edited:.2f} (Δ: {clip_score_edited - clip_score_original:.2f})', fontsize=12, fontweight='bold')
59     axes[1].axis('off')
60
61     metrics_text = f'Value of directional Similarity: {clip_directional:.2f}'
62     plt.figtext(0.5, 0.02, metrics_text, ha='center', fontsize=11, bbox=dict(boxstyle='round'))
63
64     plt.suptitle(f'Prompt: "{text_prompt}"', fontsize=13, fontweight='bold', style='italic')
65     plt.tight_layout()
66     plt.subplots_adjust(bottom=0.08)
67     plt.show()
68
69
70     result = {
71         'sample_id': f"{idx}_{prompt_idx}",
72         'sample_num': idx + 1,
73         'edit_num': prompt_idx + 1,
74         'prompt': text_prompt,
75         'clip_score_original': clip_score_original,
76         'clip_score_edited': clip_score_edited,
77         'clip_score_improvement': clip_score_edited - clip_score_original,
78         'clip_directional_similarity': clip_directional,
79     }
80
81     results.append(result)
82
83     all_edits.append({
84         'original': original_image,
85         'edited': edited_image,
86         'prompt': text_prompt,
87         'metrics': result
88     })
89
90     print(f"CLIP Score: {clip_score_edited:.2f} | improvement: {clip_score_edited - clip_score_original:.2f}")
91
92 except Exception as e:
93     print(f"Error: {e}")
94     continue
95
```


Original Image (Sample 1)

100%

50/50 [00:35<00:00, 1.39it/s]

Prompt: "change the clothes this person is wearing to blue"**Original
CLIP Score: 18.47****Edited
CLIP Score: 18.16 (Δ : -0.31)**

✓ CLIP Score: 18.16 | Improvement: -0.31 | Directional Sim: 3.76
100%
50/50 [00:36<00:00, 1.36it/s]

Prompt: "make the clothes this person is wearing red"

Original
CLIP Score: 19.26



Edited
CLIP Score: 24.88 (Δ : +5.62)



✓ CLIP Score: 24.88 | Improvement: +5.62 | Directional Sim: 16.51
100%
50/50 [00:36<00:00, 1.35it/s]

Prompt: "add stripes to the clothes"

Original
CLIP Score: 17.42



Edited
CLIP Score: 23.02 (Δ : +5.60)



✓ CLIP Score: 23.02 | Improvement: +5.60 | Directional Sim: 13.38
100%
50/50 [00:36<00:00, 1.33it/s]

Prompt: "change the clothes to a floral pattern"

Original
CLIP Score: 18.83



Edited
CLIP Score: 20.45 (Δ : +1.62)



Directional Similarity: 8.30

✓ CLIP Score: 20.45 | Improvement: +1.62 | Directional Sim: 8.30

Original Image (Sample 2)



100%

50/50 [00:37<00:00, 1.33it/s]

Prompt: "change the clothes this person is wearing to blue"

Original
CLIP Score: 22.59



Edited
CLIP Score: 26.26 (Δ : +3.67)



Directional Similarity: 3.15

✓ CLIP Score: 26.26 | Improvement: +3.67 | Directional Sim: 3.15

100%

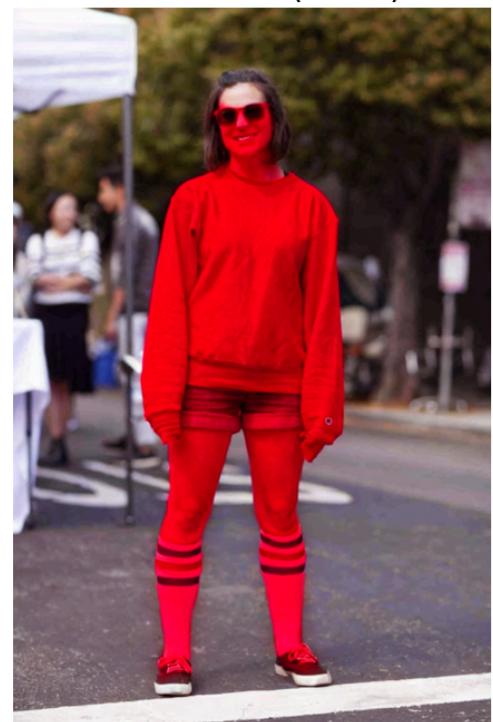
50/50 [00:37<00:00, 1.32it/s]

Prompt: "make the clothes this person is wearing red"

Original
CLIP Score: 24.07



Edited
CLIP Score: 30.63 (Δ : +6.56)



Directional Similarity: 9.25

✓ CLIP Score: 30.63 | Improvement: +6.56 | Directional Sim: 9.25

100%

50/50 [00:37<00:00, 1.32it/s]

Prompt: "add stripes to the clothes"



Directional Similarity: 6.83
 CLIP Score: 21.07 | Improvement: +2.57 | Directional Sim: 6.83
 100%
 50/50 [00:37<00:00, 1.30it/s]

Prompt: "change the clothes to a floral pattern"



Directional Similarity: 12.73
 CLIP Score: 24.86 | Improvement: +9.34 | Directional Sim: 12.73

Original Image (Sample 3)





100%

50/50 [00:37<00:00, 1.32it/s]

Prompt: "change the clothes this person is wearing to blue"

Original
CLIP Score: 20.56



Edited
CLIP Score: 19.50 (Δ : -1.06)



✓ CLIP Score: 19.50 | Improvement: -1.06 | Directional Sim: -3.27

100%

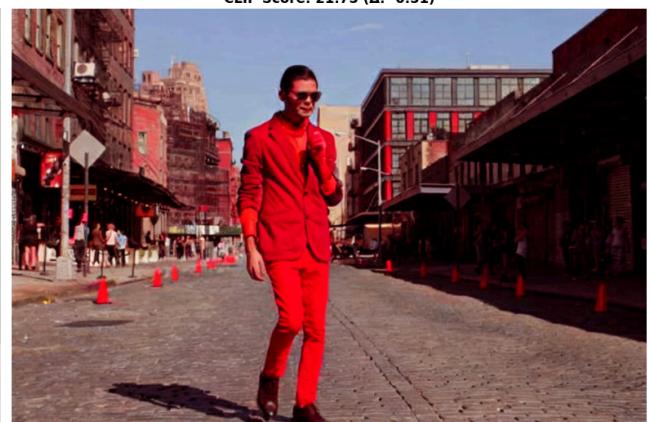
50/50 [00:37<00:00, 1.30it/s]

Prompt: "make the clothes this person is wearing red"

Original
CLIP Score: 22.26



Edited
CLIP Score: 21.75 (Δ : -0.51)





✓ CLIP Score: 23.10 | Improvement: +2.94 | Directional Sim: 7.16


```

1 lpips_value = calculate_lpips(original_image, edited_image)
2 print(f"LPIPS (baseline): {lpips_value:.4f}")

1 num_edits = len(all_edits)
2 if num_edits == 0:
3     print("Issue: no edits found in all_edits.")
4 else:
5     prompts_per_sample = 4
6     num_samples = num_edits // prompts_per_sample
7
8     fig, axes = plt.subplots(num_samples, prompts_per_sample, figsize=(14, 3 * num_samples))
9
10    if num_samples == 1:
11        axes = [axes]
12
13    for idx, edit in enumerate(all_edits):
14        row = idx // prompts_per_sample
15        col = idx % prompts_per_sample
16
17        ax = axes[row][col]
18
19        orig = edit["original"]
20        edited = edit["edited"]
21
22        small_orig = orig.copy()
23        small_edit = edited.copy()
24        small_orig.thumbnail((256, 256))
25        small_edit.thumbnail((256, 256))
26
27        combined = np.hstack((small_orig, small_edit))
28
29        ax.imshow(combined)
30        ax.set_title(edit["prompt"], fontsize=9)
31        ax.axis("off")
32
33    plt.tight_layout()
34    plt.show()
35

```

4. Fine Tuning

We first start with parameter efficient fine tuning:

```

1 gc.collect()
2 torch.cuda.empty_cache()
3
4 pipe = StableDiffusionInstructPix2PixPipeline.from_pretrained(
5     "timbrooks/instruct-pix2pix",
6     torch_dtype=torch.float16,
7     safety_checker=None,
8 )
9 pipe.to("cuda")
10
11 pipe.scheduler = DDPMscheduler.from_config(pipe.scheduler.config)
12
13 unet = pipe.unet
14 vae = pipe.vae
15 tokenizer = pipe.tokenizer
16 text_encoder = pipe.text_encoder
17 noise_scheduler = pipe.scheduler
18 device = "cuda"
19
20 for param in unet.parameters():
21     param.requires_grad = False
22
23 for name, param in unet.named_parameters():

```

```

24     should_unfreeze = False
25
26     if "attn2" in name and ("to_k" in name or "to_v" in name or "to_out" in name):
27         should_unfreeze = True
28
29     if "attn1" in name and ("to_k" in name or "to_v" in name or "to_q" in name):
30         should_unfreeze = True
31
32     if "mid_block" in name:
33         should_unfreeze = True
34
35     if "up_blocks.2" in name or "up_blocks.3" in name:
36         if "resnets" in name and "conv" in name:
37             should_unfreeze = True
38
39     if should_unfreeze:
40         param.requires_grad = True
41
42 trainable_params = [p for p in unet.parameters() if p.requires_grad]
43
44 optimizer = AdamW(trainable_params, lr=1e-5, eps=1e-4, weight_decay=0.01)
45
46 print(f"Trainable params: {sum(p.numel() for p in trainable_params):,}")
47 print(f"Total params: {sum(p.numel() for p in unet.parameters()):,}")
48 print(f"Trainable %: {100 * sum(p.numel() for p in trainable_params) / sum(p.numel() for p in unet.parameters()):.2f}%")
49 print("\nUnfrozen layers:")
50 for name, param in unet.named_parameters():
51     if param.requires_grad:
52         print(f"  {name}: {param.numel():,}")

```

[Show hidden output](#)

4.1 Training Loop for Fine-Tuning

```

1 def encode_prompt(texts):
2     tokens = tokenizer(
3         list(texts),
4         padding="max_length",
5         max_length=tokenizer.model_max_length,
6         truncation=True,
7         return_tensors="pt",
8     ).input_ids.to(device)
9     with torch.no_grad():
10         text_embeds = text_encoder(tokens)[0]
11     return text_embeds
12
13 def encode_vae_images(images):
14     images = images.to(device, dtype=torch.float16)
15     with torch.no_grad():
16         latents = vae.encode(images).latent_dist.sample()
17     return latents * vae.config.scaling_factor
18
19 @torch.no_grad()
20 def validate(val_loader, num_samples=50):
21     unet.eval()
22
23     val_losses = []
24     clip_scores = []
25     dir_sims = []
26
27     samples_processed = 0
28
29     for batch in val_loader:
30         if samples_processed >= num_samples:
31             break
32
33         pixel_values = batch["pixel_values"].to(device, dtype=torch.float16)
34         cond_images = batch["conditioning_image"].to(device, dtype=torch.float16)

```