

## Bounding-Box Experiment - Fashion Editing

The purpose of this experiment is to see how far we can push improvement in our fine-tuned fashion editing model. We want to see to what extent we can address some of the visual issues we are detecting in its edit - the most prevalent one being that it's editing outside the segmentation mask area. We want to confine the editing region stricter, so we incorporate a "hierarchical bounding" into our loss function.

The general structure of this notebook file remains the same as original, titled "Final\_Project".

### ▼ 1. Imports and Background

```

1 # imports
2 import os
3 import json
4 import torch
5 from torch import nn
6 from torch.utils.data import DataLoader
7 import torch.nn.functional as F
8 import numpy as np
9 import PIL
10 from PIL import Image
11 import matplotlib.pyplot as plt
12 from tqdm import tqdm
13 from tqdm.auto import tqdm
14 import pandas as pd
15 import requests
16 from io import BytesIO
17 import random
18
19 # models and dataset
20 from datasets import load_dataset
21 from diffusers import StableDiffusionInstructPix2PixPipeline, DPMSolverMultistepScheduler
22 from transformers import CLIPProcessor, CLIPModel
23 from datasets import load_dataset

```

```

1 # cloning pix2pix, getting it in
2
3 !git clone https://github.com/huggingface/diffusers
4
5 %cd examples/instruct_pix2pix
6 !pip install -r requirements.txt

```

[Show hidden output](#)

```

1 # from pix2pix hf website
2
3 import PIL
4 import requests
5 import torch
6 from diffusers import StableDiffusionInstructPix2PixPipeline, EulerAncestralDiscreteScheduler
7 from accelerate.utils import write_basic_config
8 write_basic_config()

```

[Show hidden output](#)

### ▼ 2. Data Preparation

```

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 #Define clip model to help us evaluate models
8 clip_model = CLIPModel.from_pretrained("openai/clip-vit-large-patch14").to("cuda")
9 clip_processor = CLIPProcessor.from_pretrained("openai/clip-vit-large-patch14")

```

[Show hidden output](#)

```

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](#)

```

1 #building triples from same person pairs
2
3 from collections import defaultdict
4
5 # Group images by person ID
6 by_pid = defaultdict(list)
7 for i in range(len(dataset["train"])):
8     by_pid[dataset["train"][i]["pid"]].append(i)
9
10 multi_outfit = {k: v for k, v in by_pid.items() if len(v) >= 2}
11 print(f"Total unique people: {len(by_pid)}")
12 print(f"People with 2+ images: {len(multi_outfit)}")
13 print(f"Potential pairs: {sum(len(v) for v in multi_outfit.values())}")

```

[Show hidden output](#)

```

1 # visualizing same person pairs
2
3 multi_outfit_pids = [pid for pid, indices in by_pid.items() if len(indices) >= 3]
4
5 fig, axes = plt.subplots(3, 3, figsize=(12, 12))
6 for row, pid in enumerate(multi_outfit_pids[:3]):
7     indices = by_pid[pid][1:3]
8     for col, idx in enumerate(indices):
9         sample = dataset["train"][idx]
10        axes[row, col].imshow(sample["images"])
11        axes[row, col].set_title(sample["caption"], fontsize=8)
12        axes[row, col].axis("off")
13
14 plt.suptitle("Same-person pairs from DeepFashion2", fontsize=14)
15 plt.tight_layout()
16 plt.show()
17

```

[Show hidden output](#)

```

1 def generate_instruction(caption1, caption2):
2     """Generate edit instruction from two captions."""
3     part2 = caption2.lower().split("in ")[-1] if "in " in caption2.lower() else caption2.lower()
4     return f"change to {part2}"
5
6 MAX_TRIPLETS = 8000
7
8 # Build triplets from same-person pairs
9 triplets = []
10 for pid, indices in tqdm(by_pid.items(), desc="Building triplets"):
11     if len(indices) < 2:
12         continue
13

```

```

14    for i in range(len(indices)):
15        for j in range(len(indices)):
16            if i == j:
17                continue
18
19            triplets.append({
20                "src_idx": indices[i],
21                "tgt_idx": indices[j],
22            })
23
24            if len(triplets) >= MAX_TRIPLETS:
25                break
26            if len(triplets) >= MAX_TRIPLETS:
27                break
28        if len(triplets) >= MAX_TRIPLETS:
29            break
30
31 print(f"Total triplets: {len(triplets)}")
32 print(f"Sample: {triplets[0]}")
33

```

[Show hidden output](#)

```

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

```

[Show hidden output](#)

```

1 import os
2 from tqdm import tqdm
3 from PIL import Image
4
5 CACHE_DIR = "deepfashion_cache"
6 os.makedirs(CACHE_DIR, exist_ok=True)
7
8 # Get unique indices we need
9 all_indices = set()
10 for t in triplets:
11     all_indices.add(t["src_idx"])
12     all_indices.add(t["tgt_idx"])
13
14 print(f"Caching {len(all_indices)} unique images...")
15
16 for idx in tqdm(all_indices, desc="Caching images"):
17     cache_path = f"{CACHE_DIR}/{idx}.png"
18     if not os.path.exists(cache_path):
19         img = dataset["train"][idx]["images"]
20         img.save(cache_path)
21

```

```
22 print("Caching complete!")
23
```

Show hidden output

```
1 from torch.utils.data import Dataset, DataLoader
2 from torchvision import transforms
3 import torchvision.transforms.functional as TF
4
5 class TripletDataset(Dataset):
6     def __init__(self, triplet_list, hf_data, size=512):
7         self.data = triplet_list
8         self.hf = hf_data
9         self.size = size
10        self.transform = transforms.Compose([
11            transforms.Resize((size, size), interpolation=transforms.InterpolationMode.BILINEAR),
12            transforms.ToTensor(),
13            transforms.Normalize([0.5, 0.5, 0.5], [0.5, 0.5, 0.5])
14        ])
15
16        self.mask_transform = transforms.Compose([
17            transforms.Resize((size, size), interpolation=transforms.InterpolationMode.NEAREST),
18            transforms.ToTensor(),
19        ])
20    # Mild color jitter shared between conditioning and target (helps generalization)
21    self.color_jitter = transforms.ColorJitter(
22        brightness=0.1,
23        contrast=0.1,
24        saturation=0.1,
25        hue=0.05,
26    )
27
28    def __len__(self):
29        return len(self.data)
30
31    def __getitem__(self, idx):
32        t = self.data[idx]
33
34        src = self.hf["train"][t["src_idx"]]
35        tgt = self.hf["train"][t["tgt_idx"]]
36
37        orig_img = src["images"]
38        edit_img = tgt["images"]
39        instruction = generate_instruction(src["caption"], tgt["caption"])
40        src_mask = src["mask"]
41
42        if orig_img.mode != "RGB":
43            orig_img = orig_img.convert("RGB")
44        if edit_img.mode != "RGB":
45            edit_img = edit_img.convert("RGB")
46        if src_mask.mode != "L":
47            src_mask = src_mask.convert("L")
48
49        if torch.rand(1).item() < 0.5:
50            orig_img = TF.hflip(orig_img)
51            edit_img = TF.hflip(edit_img)
52            src_mask = TF.hflip(src_mask)
53
54    # Mild color jitter applied consistently to both images (not the mask)
55    orig_img = self.color_jitter(orig_img)
56    edit_img = self.color_jitter(edit_img)
57
58    conditioning_image = self.transform(orig_img)
59    pixel_values = self.transform(edit_img)
60    mask_tensor = self.mask_transform(src_mask) # [1, H, W], values in [0,1]
61    mask_tensor = (mask_tensor > 0.5).float()
62
63    return {
64        "conditioning_image": conditioning_image,
65        "pixel_values": pixel_values,
```

```

66         "prompt": instruction,
67         "mask": mask_tensor,
68     }
69
70 # Create dataset
71 full_ds = TripletDataset(triplets, dataset, size=512)
72
73 train_size = int(len(full_ds) * 0.8)
74 val_size = int(len(full_ds) * 0.1)
75 test_size = len(full_ds) - train_size - val_size
76
77 train_ds, val_ds, test_ds = torch.utils.data.random_split(
78     full_ds,
79     [train_size, val_size, test_size],
80     generator=torch.Generator().manual_seed(42)
81 )
82
83 BATCH_SIZE = 1 # Small batch size
84 train_loader = DataLoader(train_ds, batch_size=BATCH_SIZE, shuffle=True, num_workers=0)
85 val_loader = DataLoader(val_ds, batch_size=BATCH_SIZE, shuffle=False, num_workers=0)
86 test_loader = DataLoader(test_ds, batch_size=BATCH_SIZE, shuffle=False, num_workers=0)
87
88 print(f"Train: {len(train_ds)}, Val: {len(val_ds)}, Test: {len(test_ds)}")
89
90 test_batch = next(iter(train_loader))
91 print(f"Conditioning shape: {test_batch['conditioning_image'].shape}")
92 print(f"Sample prompt: {test_batch['prompt'][0]}")
93 print("DataLoader working!")

```

Show hidden output

## ▼ 3. Model Preparation and Baseline

### ▼ 3.1 Calculating Clip Score and Directional Similarity Scores (Evaluation 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

```

## 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     results.append(result)
81
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

```

[Show hidden output](#)

## 4. 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 Loops (Edited to include Bounding boxes and zoning from Fashionpedia)

```

1 from torch.optim.lr_scheduler import CosineAnnealingLR
2 from torchvision.transforms.functional import to_pil_image
3
4 def encode_prompt(texts):
5     tokens = tokenizer(
6         list(texts),
7         padding="max_length",
8         max_length=tokenizer.model_max_length,
9         truncation=True,
10        return_tensors="pt",
11    ).input_ids.to(device)
12    with torch.no_grad():
13        text_embeds = text_encoder(tokens)[0]
14    return text_embeds
15
16 def encode_vae_images(images):
17    images = images.to(device, dtype=torch.float16)
18    with torch.no_grad():
19        latents = vae.encode(images).latent_dist.sample()
20    return latents * vae.config.scaling_factor
21
22 def create_bbox_mask(bboxes, image_size, expansion_factor=1.15):
23     """
24     Author: Disclaimer! I Used ChatGPT to help me generate this function. The
25     source idea and inspiration came from Bayesian Approach to Digital Matting

```

```

26     (Chuang Et. al 2001)
27
28     Create soft masks from bounding boxes with slight expansion for context.
29
30     Args:
31         bboxes: List of [x1, y1, x2, y2] bounding boxes (normalized 0-1 or pixel coords)
32         image_size: (height, width) of the image
33         expansion_factor: Factor to expand bbox (1.15 = 15% larger)
34     """
35     h, w = image_size
36     batch_size = len(bboxes)
37     masks = torch.zeros(batch_size, 1, h, w, device=device, dtype=torch.float16)
38
39     for i, bbox in enumerate(bboxes):
40         if bbox is None or len(bbox) != 4:
41             continue
42
43         x1, y1, x2, y2 = bbox
44         # If normalized coordinates, convert to pixels
45         if x2 <= 1.0:
46             x1, x2 = x1 * w, x2 * w
47             y1, y2 = y1 * h, y2 * h
48         # Expand bbox slightly
49         cx, cy = (x1 + x2) / 2, (y1 + y2) / 2
50         bw, bh = (x2 - x1) * expansion_factor, (y2 - y1) * expansion_factor
51         x1_new = max(0, int(cx - bw/2))
52         y1_new = max(0, int(cy - bh/2))
53         x2_new = min(w, int(cx + bw/2))
54         y2_new = min(h, int(cy + bh/2))
55         masks[i, 0, y1_new:y2_new, x1_new:x2_new] = 1.0
56
57     return masks
58
59 def compute_hierarchical_loss(noise_pred, noise, mask_latent, bbox_latent,
60                               lambda_context=0.3, lambda_preserve=0.5):
61     """
62     Author: Disclaimer! I Used ChatGPT to help me generate this function. The
63     source idea and inspiration came from Bayesian Approach to Digital Matting
64     (Chuang Et. al 2001)
65
66     Three-zone hierarchical loss:
67     - Tight clothing region (from segmentation): full editing loss
68     - Context region: moderate preservation (allows smooth transitions)
69     - Background: strong preservation
70     """
71     # Zone 1: Tight clothing mask (primary edit region)
72     M_tight = mask_latent
73
74     # Zone 2: Transition Zone
75     M_context = bbox_latent * (1.0 - mask_latent)
76
77     # Zone 3: Background
78     M_bg = 1.0 - bbox_latent
79
80     # Compute losses for each zone
81     L_tight = F.mse_loss((noise_pred * M_tight).float(), (noise * M_tight).float())
82     L_context = F.mse_loss((noise_pred * M_context).float(), (noise * M_context).float())
83     L_bg = F.mse_loss((noise_pred * M_bg).float(), (noise * M_bg).float())
84
85     # Weighted combination
86     total_loss = L_tight + lambda_context * L_context + lambda_preserve * L_bg
87
88     return total_loss, {
89         'L_tight': L_tight.item(),
90         'L_context': L_context.item(),
91         'L_bg': L_bg.item()
92     }
93
94 @torch.no_grad()
95 def validate(val_loader, num_samples=50):

```

```

96     """Calculate validation metrics.
97         Disclaimer: I used ChatGPT to help me generate part of this function
98     """
99     unet.eval()
100
101    val_losses = []
102    clip_scores = []
103    dir_sims = []
104
105    samples_processed = 0
106    for batch in val_loader:
107        if samples_processed >= num_samples: #no samples
108            break
109        pixel_values = batch["pixel_values"].to(device, dtype=torch.float16)
110        cond_images = batch["conditioning_image"].to(device, dtype=torch.float16)
111        masks = batch["mask"].to(device, dtype=torch.float16)
112        prompts = batch["prompt"]
113        bboxes = batch.get("bbox", None)
114
115        # Compute validation loss
116        target_latents = encode_vae_images(pixel_values)
117        cond_latents = encode_vae_images(cond_images)
118        mask_latent = F.interpolate(masks, size=target_latents.shape[-2:], mode="nearest")
119        mask_latent = mask_latent.expand(-1, 4, -1, -1)
120
121        # Create bounding boxes if available
122        if bboxes is not None:
123            bbox_masks = create_bbox_mask(bboxes, (masks.shape[2], masks.shape[3]))
124            bbox_latent = F.interpolate(bbox_masks, size=target_latents.shape[-2:], mode="nearest")
125            bbox_latent = bbox_latent.expand(-1, 4, -1, -1)
126        else:
127            # I used CHATGPT to help me design this fallback
128            bbox_latent = F.max_pool2d(mask_latent, kernel_size=5, stride=1, padding=2)
129
130        noise = torch.randn_like(target_latents)
131        timesteps = torch.randint(0, noise_scheduler.config.num_train_timesteps,
132                                  (target_latents.shape[0],), device=device, dtype=torch.long)
133        noisy_latents = noise_scheduler.add_noise(target_latents, noise, timesteps)
134        model_input = torch.cat([noisy_latents, cond_latents], dim=1)
135        text_embeds = encode_prompt(prompts)
136
137        with torch.autocast(device_type="cuda", dtype=torch.float16):
138
139            #calculate the loss
140            noise_pred = unet(model_input, timesteps, encoder_hidden_states=text_embeds, return_dict=False)
141            loss, _ = compute_hierarchical_loss(noise_pred, noise, mask_latent, bbox_latent,
142                                                lambda_context=LAMBDA_CONTEXT,
143                                                lambda_preserve=LAMBDA_PRESERVE)
144            val_losses.append(loss.item())
145            if samples_processed % 5 == 0: #skip every 5 to save compute time
146                cond_img_pil = to_pil_image(cond_images[0].cpu().float() * 0.5 + 0.5)
147
148                with torch.autocast("cuda"):
149                    edited = pipe(
150                        prompt=prompts[0],
151                        image=cond_img_pil,
152                        num_inference_steps=20,
153                        guidance_scale=7.5,
154                        image_guidance_scale=1.5,
155                    ).images[0]
156
157                clip_score = calculate_clip_score(edited, prompts[0], clip_model, clip_processor)
158                dir_sim = calculate_clip_directional_similarity(cond_img_pil, edited, prompts[0], clip_processor)
159
160                clip_scores.append(clip_score)
161                dir_sims.append(dir_sim)
162
163                samples_processed += 1
164
165        unet.train()

```

```

166
167     return {
168         "val_loss": sum(val_losses) / len(val_losses),
169         "clip_score": sum(clip_scores) / len(clip_scores) if clip_scores else 0,
170         "dir_sim": sum(dir_sims) / len(dir_sims) if dir_sims else 0,
171     }
172
173
174 # Hyperparameters
175 EPOCHS = 4
176 MAX_GRAD_NORM = 1.0
177 LAMBDA_CONTEXT = 0.5 # Weight for context region
178 LAMBDA_PRESERVE = 0.7 # Weight for background preservation
179 PATIENCE = 2
180
181 num_training_steps = EPOCHS * len(train_loader)
182 scheduler = CosineAnnealingLR(optimizer, T_max=num_training_steps)
183
184 unet.train()
185 train_losses = []
186 val_metrics_history = []
187 best_val_dirsim = float("-inf")
188 epochs_no_improve = 0
189
190 for epoch in range(EPOCHS):
191     epoch_losses = []
192     loss_components = {'tight': [], 'context': [], 'bg': []}
193     pbar = tqdm(train_loader, desc=f"Epoch {epoch+1}/{EPOCHS}")
194
195     for step, batch in enumerate(pbar):
196         pixel_values = batch["pixel_values"].to(device, dtype=torch.float16)
197         cond_images = batch["conditioning_image"].to(device, dtype=torch.float16)
198         masks = batch["mask"].to(device, dtype=torch.float16)
199         prompts = batch["prompt"]
200
201         bboxes = batch.get("bbox", None)
202
203         with torch.no_grad():
204             target_latents = encode_vae_images(pixel_values)
205             cond_latents = encode_vae_images(cond_images)
206             mask_latent = F.interpolate(masks, size=target_latents.shape[-2:], mode="nearest")
207             mask_latent = mask_latent.expand(-1, 4, -1, -1)
208             if bboxes is not None:
209                 bbox_masks = create_bbox_mask(bboxes, (masks.shape[2], masks.shape[3]))
210                 bbox_latent = F.interpolate(bbox_masks, size=target_latents.shape[-2:], mode="nearest")
211                 bbox_latent = bbox_latent.expand(-1, 4, -1, -1)
212             else:
213                 # Fallback: use dilated segmentation mask as approximation from chatGPT
214                 bbox_latent = F.max_pool2d(mask_latent, kernel_size=5, stride=1, padding=2)
215
216             noise = torch.randn_like(target_latents)
217             bsz = target_latents.shape[0]
218             timesteps = torch.randint(
219                 0, noise_scheduler.config.num_train_timesteps,
220                 (bsz,), device=device, dtype=torch.long
221             )
222
223             noisy_latents = noise_scheduler.add_noise(target_latents, noise, timesteps)
224             model_input = torch.cat([noisy_latents, cond_latents], dim=1)
225             text_embeds = encode_prompt(prompts)
226
227             optimizer.zero_grad()
228             with torch.autocast(device_type="cuda", dtype=torch.float16):
229                 noise_pred = unet(
230                     model_input,
231                     timesteps,
232                     encoder_hidden_states=text_embeds,
233                     return_dict=False
234                 )[0]
235                 loss, loss_dict = compute_hierarchical_loss(

```

```

236         noise_pred, noise, mask_latent, bbox_latent,
237         lambda_context=LAMBDA_CONTEXT,
238         lambda_preserve=LAMBDA_PRESERVE
239     )
240
241
242     loss_components['tight'].append(loss_dict['L_tight'])
243     loss_components['context'].append(loss_dict['L_context'])
244     loss_components['bg'].append(loss_dict['L_bg'])
245     if not torch.isfinite(loss):
246         print(f"Non-finite loss at step {step}, skipping...")
247         continue
248
249     loss.backward()
250     torch.nn.utils.clip_grad_norm_(trainable_params, MAX_GRAD_NORM)
251
252     # Check for NaN gradients
253     has_nan_grad = False
254     for p in trainable_params:
255         if p.grad is not None and not torch.isfinite(p.grad).all():
256             has_nan_grad = True
257             break
258
259     if has_nan_grad:
260         print(f"Non-finite gradient at step {step}, skipping...")
261         optimizer.zero_grad()
262         continue
263
264     optimizer.step()
265     scheduler.step()
266
267     epoch_losses.append(loss.item())
268     pbar.set_postfix({
269         "loss": f"{loss.item():.4f}",
270         "L_tight": f"{loss_dict['L_tight']:.3f}",
271         "L_ctx": f"{loss_dict['L_context']:.3f}",
272         "lr": f"{scheduler.get_last_lr()[0]:.2e}"
273     })
274
275     avg_loss = sum(epoch_losses) / len(epoch_losses) if epoch_losses else float('nan')
276     train_losses.append(avg_loss)
277
278     # ChatGPT generated these print statements for me
279     print(f"Epoch {epoch+1} average loss: {avg_loss:.4f}")
280     print(f" Loss breakdown - Tight: {sum(loss_components['tight'])/len(loss_components['tight'])}")
281     print(f"Context: {sum(loss_components['context'])/len(loss_components['context']):.4f}, ")
282     print(f"Background: {sum(loss_components['bg'])/len(loss_components['bg']):.4f}")
283
284     print(f"Running validation...")
285     val_metrics = validate(val_loader, num_samples=50)
286     val_metrics_history.append(val_metrics)
287
288     print(f" Val Loss: {val_metrics['val_loss']:.4f}")
289     print(f" CLIP Score: {val_metrics['clip_score']:.2f}")
290     print(f" Directional Sim: {val_metrics['dir_sim']:.2f}")
291
292
293 #Early stopping if no improvement
294 if val_metrics['dir_sim'] > best_val_dirsim:
295     best_val_dirsim = val_metrics['dir_sim']
296     epochs_no_improve = 0
297     torch.save(unet.state_dict(), "unet_best.pt")
298     print(f" New best model saved! (dir_sim: {best_val_dirsim:.2f})")
299 else:
300     epochs_no_improve += 1
301     print(f" No improvement for {epochs_no_improve} epoch(s)")
302
303     if epochs_no_improve >= PATIENCE:
304         print(f"Early stopping triggered at epoch {epoch+1}")
305         break

```

```

306
307 #save at this checkpoint
308 torch.save({
309     'epoch': epoch,
310     'unet_state_dict': unet.state_dict(),
311     'optimizer_state_dict': optimizer.state_dict(),
312     'scheduler_state_dict': scheduler.state_dict(),
313     'train_losses': train_losses,
314     'val_metrics': val_metrics_history,
315 }, f"checkpoint_epoch{epoch+1}.pt")
316
317 print("Training complete!")
318 print(f"Best directional similarity: {best_val_dirsim:.2f}")

```

Show hidden output

```

1 # plot training loss
2
3 plt.figure(figsize=(10, 5))
4 plt.plot(range(1, len(train_losses) + 1), train_losses, marker='o')
5 plt.xlabel("Epoch")
6 plt.ylabel("Average Loss")
7 plt.title("Training Loss - Bounding Box Fine-tuning")
8 plt.grid(True)
9 plt.show()

```

```

NameError Traceback (most recent call last)
/tmp/ipython-input-2932974737.py in <cell line: 0>()
      1 # plot training loss
      2
----> 3 plt.figure(figsize=(10, 5))
      4 plt.plot(range(1, len(train_losses) + 1), train_losses, marker='o')
      5 plt.xlabel("Epoch")

NameError: name 'plt' is not defined

```

```

1 # Load fine-tuned model for eval
2
3 state_dict = torch.load("unet_best.pt", map_location=device)
4 unet.load_state_dict(state_dict, strict=False)
5
6 pipe.unet = unet
7 pipe.to("cuda")
8 pipe.unet.eval()
9
10 print("Fine-tuned model loaded!")

```

Fine-tuned model loaded!

```

1 custom_prompts = [
2     "change to a floral pattern",
3     "make the shirt red",
4     "add stripes to the clothes",
5     "change shirt to blue",
6 ]
7
8 NUM_TEST_IMAGES = 5
9
10 fig, axes = plt.subplots(NUM_TEST_IMAGES, len(custom_prompts) + 1, figsize=(18, 4 * NUM_TEST_IMAGES))
11
12 custom_results = []
13
14 for row in range(NUM_TEST_IMAGES):
15     actual_idx = test_ds.indices[row]
16     t = triplets[actual_idx]
17     src = dataset["train"][t["src_idx"]]
18     original_image = src["images"]
19
20     if max(original_image.size) > 512:

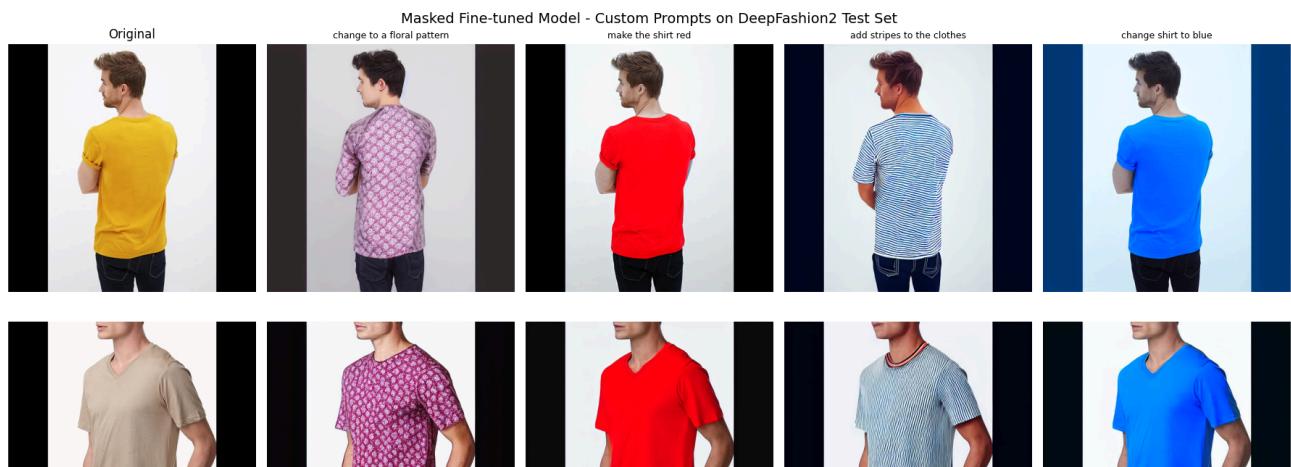
```

```
21     original_image = original_image.copy()
22     original_image.thumbnail((512, 512), Image.Resampling.LANCZOS)
23
24     axes[row, 0].imshow(original_image)
25     axes[row, 0].set_title("Original" if row == 0 else "")
26     axes[row, 0].axis("off")
27
28     for col, prompt in enumerate(custom_prompts):
29         with torch.autocast("cuda"):
30             edited = pipe(
31                 prompt=prompt,
32                 image=original_image,
33                 num_inference_steps=50,
34                 guidance_scale=7.5,
35                 image_guidance_scale=1.5,
36             ).images[0]
37
38     clip_score = calculate_clip_score(edited, prompt, clip_model, clip_processor)
39     clip_directional = calculate_clip_directional_similarity(
40         original_image, edited, prompt, clip_model, clip_processor
41     )
42     lpips_val = calculate_lpips(original_image, edited)
43
44     custom_results.append({
45         'image': row,
46         'prompt': prompt,
47         'clip_score': clip_score,
48         'directional_similarity': clip_directional,
49         'lpips': lpips_val,
50     })
51
52     axes[row, col + 1].imshow(edited)
53     axes[row, col + 1].set_title(prompt if row == 0 else "", fontsize=9)
54     axes[row, col + 1].axis("off")
55
56 plt.tight_layout()
57 plt.savefig("custom_prompts_deepfashion2.png", dpi=150, bbox_inches='tight')
58 plt.show()
59
60 # Summary
61 df_custom = pd.DataFrame(custom_results)
62 print("\n Custom Prompts Evaluation ===")
63 print(f"mean CLIP Score: {df_custom['clip_score'].mean():.2f}")
64 print(f"mean Directional Similarity: {df_custom['directional_similarity'].mean():.2f}")
65 print(f"mean LPIPS: {df_custom['lpips'].mean():.4f}")
66 print("\nPer-prompt averages:")
67 print(df_custom.groupby('prompt')[['clip_score', 'directional_similarity']].mean())
```

```

100% 50/50 [00:11<00:00, 4.33it/s]
100% 50/50 [00:11<00:00, 4.24it/s]
100% 50/50 [00:11<00:00, 4.24it/s]
100% 50/50 [00:11<00:00, 4.30it/s]
100% 50/50 [00:11<00:00, 4.36it/s]
100% 50/50 [00:11<00:00, 4.38it/s]
100% 50/50 [00:11<00:00, 4.40it/s]
100% 50/50 [00:11<00:00, 4.36it/s]
100% 50/50 [00:11<00:00, 4.35it/s]
100% 50/50 [00:11<00:00, 4.34it/s]
100% 50/50 [00:11<00:00, 4.32it/s]
100% 50/50 [00:11<00:00, 4.33it/s]
100% 50/50 [00:11<00:00, 4.34it/s]
100% 50/50 [00:11<00:00, 4.36it/s]
100% 50/50 [00:11<00:00, 4.36it/s]
100% 50/50 [00:11<00:00, 4.35it/s]
100% 50/50 [00:11<00:00, 4.33it/s]

```



```

1 # compare baseline vs fine tuned
2 print(f"Baseline Mean CLIP Score: {pd.DataFrame(results)['clip_score_edited'].mean():.2f}")
3 print(f"Fine-tuned Mean CLIP Score: {df_custom['clip_score'].mean():.2f}")
4 print(f"\nBaseline Mean Directional Sim: {pd.DataFrame(results)['clip_directional_similarity'].mean():.2f}")
5 print(f"Fine-tuned Mean Directional Sim: {df_custom['directional_similarity'].mean():.2f}")

```



```

1 # evaluating on fashion pedia dataset
2
3 test_prompts = [
4     "change to a floral pattern",
5     "make the shirt red",
6     "add stripes to the clothes",
7     "change shirt to blue",
8 ]
9

```

```
10 NUM_SAMPLES = 5
11
12 finetuned_fashionpedia_results = []
13
14 for idx in range(NUM_SAMPLES):
15     sample = fashionpedia_samples[idx]
16     original_image = sample['image']
17
18     if max(original_image.size) > 512:
19         original_image = original_image.copy()
20         original_image.thumbnail((512, 512), Image.Resampling.LANCZOS)
21
22     for prompt in test_prompts:
23         # Generate with fine-tuned model
24         with torch.autocast("cuda"):
25             edited = pipe(
26                 prompt=prompt,
27                 image=original_image,
28                 num_inference_steps=50,
29                 guidance_scale=7.5,
30                 image_guidance_scale=1.5,
31             ).images[0]
32
33     # Calculate metrics
34     clip_score = calculate_clip_score(edited, prompt, clip_model, clip_processor)
35     clip_directional = calculate_clip_directional_similarity(
36         original_image, edited, prompt, clip_model, clip_processor
37     )
38
39     finetuned_fashionpedia_results.append({
40         'sample': idx,
41         'prompt': prompt,
42         'clip_score': clip_score,
43         'directional_similarity': clip_directional,
44     })
45
46 df_finetuned_fp = pd.DataFrame(finetuned_fashionpedia_results)
47
48 print(f"Mean CLIP Score: {df_finetuned_fp['clip_score'].mean():.2f}")
49 print(f"Mean Directional Similarity: {df_finetuned_fp['directional_similarity'].mean():.2f}")
50
51 # Compare with baseline on FashionPedia
52 print(f"Baseline Mean CLIP Score: {pd.DataFrame(results)['clip_score_edited'].mean():.2f}")
53 print(f"Fine-tuned Mean CLIP Score: {df_finetuned_fp['clip_score'].mean():.2f}")
54 print(f"\nBaseline Mean Directional Sim: {pd.DataFrame(results)['clip_directional_similarity'].mean():.2f}")
55 print(f"Fine-tuned Mean Directional Sim: {df_finetuned_fp['directional_similarity'].mean():.2f}")
```

```

100%                                         50/50 [00:08<00:00, 6.30it/s]
100%                                         50/50 [00:08<00:00, 6.24it/s]
100%                                         50/50 [00:08<00:00, 6.23it/s]
100%                                         50/50 [00:08<00:00, 6.24it/s]
100%                                         50/50 [00:08<00:00, 6.29it/s]
100%                                         50/50 [00:08<00:00, 6.29it/s]
100%                                         50/50 [00:08<00:00, 6.28it/s]
100%                                         50/50 [00:08<00:00, 6.27it/s]
100%                                         50/50 [00:08<00:00, 6.28it/s]
100%                                         50/50 [00:08<00:00, 6.29it/s]
100%                                         50/50 [00:08<00:00, 6.29it/s]
100%                                         50/50 [00:08<00:00, 6.29it/s]
100%                                         50/50 [00:08<00:00, 6.26it/s]
100%                                         50/50 [00:08<00:00, 6.25it/s]
100%                                         50/50 [00:08<00:00, 6.32it/s]
100%                                         50/50 [00:08<00:00, 6.28it/s]
100%                                         50/50 [00:09<00:00, 5.04it/s]
100%                                         50/50 [00:10<00:00, 5.04it/s]
100%                                         50/50 [00:10<00:00, 5.03it/s]
100%                                         50/50 [00:10<00:00, 5.03it/s]

==== Fine-tuned Model on FashionPedia ====
Mean CLIP Score: 19.16
Mean Directional Similarity: 10.61

==== Comparison on FashionPedia ====
Baseline Mean CLIP Score: 23.03
Fine-tuned Mean CLIP Score: 19.16

Baseline Mean Directional Sim: 7.82
Fine-tuned Mean Directional Sim: 10.61

```

```

1 # plotting for fine tuned model on fashionpedia
2
3 fig, axes = plt.subplots(NUM_SAMPLES, len(test_prompts) + 1, figsize=(18, 4 * NUM_SAMPLES))
4
5 for row in range(NUM_SAMPLES):
6     sample = fashionpedia_samples[row]
7     original_image = sample['image']
8
9     if max(original_image.size) > 512:
10        original_image = original_image.copy()
11        original_image.thumbnail((512, 512), Image.Resampling.LANCZOS)
12
13    axes[row, 0].imshow(original_image)
14    axes[row, 0].set_title("Original" if row == 0 else "")
15    axes[row, 0].axis("off")
16
17    for col, prompt in enumerate(test_prompts):
18        with torch.autocast("cuda"):
19            edited = pipe(
20                prompt=prompt,
21                image=original_image,
22                num_inference_steps=50,
23                guidance_scale=7.5,
24                image_guidance_scale=1.5,
25            ).images[0]
26
27        axes[row, col + 1].imshow(edited)
28        axes[row, col + 1].set_title(prompt if row == 0 else "", fontsize=9)
29        axes[row, col + 1].axis("off")

```

```

30
31 plt.suptitle("Fine-tuned Model on FashionPedia (unseen data)", fontsize=14)
32 plt.tight_layout()
33 plt.savefig("finetuned_on_fashionpedia.png", dpi=150, bbox_inches='tight')
34 plt.show()

```

100%	50/50 [00:08<00:00, 6.28it/s]
100%	50/50 [00:08<00:00, 6.24it/s]
100%	50/50 [00:08<00:00, 6.27it/s]
100%	50/50 [00:08<00:00, 6.30it/s]
100%	50/50 [00:08<00:00, 6.28it/s]
100%	50/50 [00:08<00:00, 6.25it/s]
100%	50/50 [00:08<00:00, 6.23it/s]
100%	50/50 [00:08<00:00, 6.25it/s]
100%	50/50 [00:08<00:00, 6.25it/s]
100%	50/50 [00:08<00:00, 6.24it/s]
100%	50/50 [00:08<00:00, 6.28it/s]
100%	50/50 [00:08<00:00, 6.25it/s]
100%	50/50 [00:08<00:00, 6.25it/s]
100%	50/50 [00:08<00:00, 6.28it/s]
100%	50/50 [00:08<00:00, 6.25it/s]
100%	50/50 [00:08<00:00, 6.26it/s]
100%	50/50 [00:10<00:00, 5.03it/s]
100%	50/50 [00:10<00:00, 5.05it/s]
100%	50/50 [00:10<00:00, 5.02it/s]
100%	50/50 [00:10<00:00, 5.04it/s]



5. The following code is used to generate some visualizations for our paper given the datasets we were using (FashionPedia). It is mostly generated using AI and not use in our actual experimentation.



```

1 import matplotlib.pyplot as plt
2 import matplotlib.patches as patches
3 from datasets import load_dataset
4
5 # 1. Load dataset
6 print("Loading FashionPedia stream...")
7 dataset_fashionpedia = load_dataset("detection-datasets/fashionpedia", split="train", streaming=True)
8

```

```
9 # 2. Grab just 5 samples
10 print("Fetching samples...")
11 fashionpedia_samples = list(dataset_fashionpedia.take(5))
12
13 # 3. Robust Helper to get Class Names
14 # We check if 'objects' is a dict or an object to avoid the AttributeError
15 def get_category_names(dataset):
16     try:
17         # Try standard Hugging Face object access
18         return dataset.features['objects'].feature['category'].names
19     except AttributeError:
20         try:
21             # Fallback: Try dictionary access (common in some Colab envs)
22             return dataset.features['objects'][feature]['category']['names']
23         except (KeyError, TypeError):
24             # If all else fails, return None and we will just print IDs
25             print("Warning: Could not extract class names. Displaying IDs only.")
26             return None
27
28 class_names = get_category_names(dataset_fashionpedia)
29
30 def get_label_name(cat_id):
31     if class_names and cat_id < len(class_names):
32         return class_names[cat_id]
33     return str(cat_id) # Return ID if names failed to load
34
35 # 4. Visualization Function
36 def visualize_sample(sample):
37     image = sample['image']
38     objects = sample['objects']
39
40     # Create a figure
41     fig, ax = plt.subplots(figsize=(10, 10))
42     ax.imshow(image)
43
44     # Loop through all objects detected in the image
45     # Note: Some datasets use 'categories' instead of 'category', we handle both
46     cat_keys = objects.get('category', objects.get('categories', []))
47
48     for bbox, cat_id in zip(objects['bbox'], cat_keys):
49         # Format is [x_min, y_min, width, height]
50         x, y, w, h = bbox
51
52         # Create a Rectangle patch
53         rect = patches.Rectangle((x, y), w, h, linewidth=2, edgecolor="#00FF00", facecolor='none')
54         ax.add_patch(rect)
55
56         # Add the label text above the box
57         label_text = get_label_name(cat_id)
58         ax.text(x, y, label_text, color='white', fontsize=9, weight='bold',
59                 bbox=dict(facecolor='#00FF00', edgecolor='none', alpha=0.5))
60
61     plt.axis('off')
62     plt.show()
63
64 # 5. Show the first image
65 print(f"Visualizing Sample 0...")
66 visualize_sample(fashionpedia_samples[0])
```

Loading FashionPedia stream...  
 Fetching samples...  
 Warning: Could not extract class names. Displaying IDs only.  
 Visualizing Sample 0...



```

1 import matplotlib.pyplot as plt
2 import matplotlib.patches as patches
3 import numpy as np
4 from datasets import load_dataset
5
6 # 1. Load Dataset
7 print("Loading FashionPedia stream...")
8 dataset = load_dataset("detection-datasets/fashionpedia", split="train", streaming=True)
9 sample = next(iter(dataset))
10
11 image_pil = sample['image']
12 image = np.array(image_pil)
13 h, w, _ = image.shape
14 objects = sample['objects']
15
16 # get the dress (target_idx 3)
17 max_area = 0
18 target_idx = 3
19 bx, by, bw, bh = [int(v) for v in objects['bbox'][target_idx]]
20
21 # 3. Create Zone Masks
22
23 # Zone 1: Tight Edit Region (the actual bounding box)
24 mask_tight = np.zeros((h, w), dtype=np.uint8)
25 mask_tight[by:by+bh, bx:bx+bw] = 1
26
27 # Zone 2: Context Region (expanded box minus tight box)

```

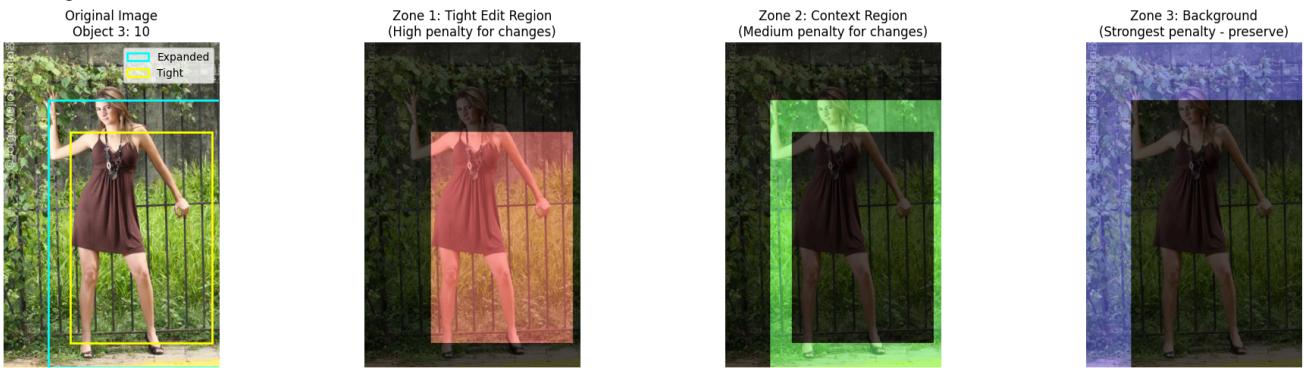
```

28 EXPANSION_FACTOR = 1.3
29 cx, cy = bx + bw/2, by + bh/2
30 new_w, new_h = bw * EXPANSION_FACTOR, bh * EXPANSION_FACTOR
31
32 nx1 = int(max(0, cx - new_w/2))
33 ny1 = int(max(0, cy - new_h/2))
34 nx2 = int(min(w, cx + new_w/2))
35 ny2 = int(min(h, cy + new_h/2))
36
37 mask_expanded = np.zeros((h, w), dtype=np.uint8)
38 mask_expanded[ny1:ny2, nx1:nx2] = 1
39
40 mask_context = mask_expanded - mask_tight
41 mask_context = np.clip(mask_context, 0, 1)
42
43 # Zone 3: Background (everything outside expanded box)
44 mask_bg = 1 - mask_expanded
45
46 # 4. Visualization
47 fig, axes = plt.subplots(1, 4, figsize=(20, 5))
48
49 # Plot A: Original Image with bounding boxes
50 axes[0].imshow(image)
51 # Original bbox (yellow)
52 rect_tight = patches.Rectangle((bx, by), bw, bh, linewidth=2,
53                                 edgecolor='yellow', facecolor='none', label='Tight')
54 # Expanded bbox (cyan)
55 rect_expanded = patches.Rectangle((nx1, ny1), nx2-nx1, ny2-ny1, linewidth=2,
56                                   edgecolor='cyan', facecolor='none', label='Expanded')
57 axes[0].add_patch(rect_expanded)
58 axes[0].add_patch(rect_tight)
59 axes[0].set_title(f"Original Image\nObject {target_idx}: {objects['category'][target_idx]}")
60 axes[0].legend()
61 axes[0].axis('off')
62
63 # Plot B: Zone 1 (Tight Edit Region)
64 overlay_tight = np.zeros_like(image)
65 overlay_tight[mask_tight == 1] = [255, 100, 100] # Red tint
66 blended = np.where(mask_tight[...], None) == 1,
67             image * 0.5 + overlay_tight * 0.5,
68             image * 0.3).astype(np.uint8)
69 axes[1].imshow(blended)
70 axes[1].set_title("Zone 1: Tight Edit Region\n(High penalty for changes)")
71 axes[1].axis('off')
72
73 # Plot C: Zone 2 (Context Region)
74 overlay_context = np.zeros_like(image)
75 overlay_context[mask_context == 1] = [100, 255, 100] # Green tint
76 blended = np.where(mask_context[...], None) == 1,
77             image * 0.5 + overlay_context * 0.5,
78             image * 0.3).astype(np.uint8)
79 axes[2].imshow(blended)
80 axes[2].set_title("Zone 2: Context Region\n(Medium penalty for changes)")
81 axes[2].axis('off')
82
83 # Plot D: Zone 3 (Background - Should be Preserved)
84 overlay_bg = np.zeros_like(image)
85 overlay_bg[mask_bg == 1] = [100, 100, 255] # Blue tint
86 blended = np.where(mask_bg[...], None) == 1,
87             image * 0.5 + overlay_bg * 0.5,
88             image * 0.3).astype(np.uint8)
89 axes[3].imshow(blended)
90 axes[3].set_title("Zone 3: Background\n(Strongest penalty - preserve)")
91 axes[3].axis('off')
92
93 plt.tight_layout()
94 plt.show()
95
96 # Print zone statistics
97 print(f"\nZone Statistics:")

```

```
98 print(f"Zone 1 (Tight): {mask_tight.sum()} pixels ({100*mask_tight.sum()/(h*w):.1f}%)")
99 print(f"Zone 2 (Context): {mask_context.sum()} pixels ({100*mask_context.sum()/(h*w):.1f}%)"
100 print(f"Zone 3 (Background): {mask_bg.sum()} pixels ({100*mask_bg.sum()/(h*w):.1f}%)")
```

Loading FashionPedia stream...



Zone Statistics:

```
Zone 1 (Tight): 297920 pixels (42.7%)
Zone 2 (Context): 156760 pixels (22.4%)
Zone 3 (Background): 243688 pixels (34.9%)
```

```
1 # --- Load InstructPix2Pix ---
2 from diffusers import StableDiffusionInstructPix2PixPipeline, DPMSolverMultistepScheduler
3 from datasets import load_dataset
4 import torch
5 import PIL
6 import matplotlib.pyplot as plt
7
8 # Load model
9 model_id = "timbrooks/instruct-pix2pix"
10 pipe = StableDiffusionInstructPix2PixPipeline.from_pretrained(
11     model_id,
12     torch_dtype=torch.float16,
13     safety_checker=None,
14 )
15 pipe.to("cuda")
16
17 pipe.scheduler = DPMSolverMultistepScheduler.from_config(pipe.scheduler.config)
18
19 # --- Load ONE example from FashionPedia ---
20 dataset = load_dataset("detection-datasets/fashionpedia", split="train", streaming=True)
21 first_item = next(iter(dataset))      # get only the FIRST sample
22
23 original_image = first_item["image"].convert("RGB")
24
25 print("Loaded first FashionPedia image")
26 display(original_image)
27
28 # --- Apply edit ---
29 prompt = "Give her clothes a floral pattern"
30
31 edited = pipe(
32     prompt=prompt,
33     image=original_image,
34     num_inference_steps=100,
35     guidance_scale=12,
36     image_guidance_scale=2.0,
```

```
37 ).images[0]
38
39 # --- Show side-by-side ---
40 plt.figure(figsize=(12,6))
41 plt.subplot(1,2,1)
42 plt.title("Original")
43 plt.imshow(original_image)
44 plt.axis("off")
45
46 plt.subplot(1,2,2)
47 plt.title("Edited: Floral Pattern")
48 plt.imshow(edited)
49 plt.axis("off")
50
51 plt.show()
52
53
```