

✓ Ablation Study for Domain-Specific InstructPix2Pix

This standalone notebook reuses the *Final_Project* pipeline (data prep, etc), but runs the ablation study.

```
1 # install packages required
2 %pip install -q lpips pytorch-msssim
```

```
1 import os
2 import gc
3 import math
4 import json
5 import copy
6 import random
7 import re
8 from dataclasses import dataclass, asdict
9 from typing import Dict, List, Any
10
11 import numpy as np
12 import pandas as pd
13 import torch
14 import torch.nn.functional as F
15 from torch import nn
16 from torch.utils.data import Dataset, DataLoader
17 from torchvision import transforms
18 import torchvision.transforms.functional as TF
19 from torchvision.transforms.functional import to_pil_image
20
21 from PIL import Image
22 import matplotlib.pyplot as plt
23 from tqdm.auto import tqdm
24 from collections import defaultdict
25
26 from datasets import load_dataset
27 from diffusers import StableDiffusionInstructPix2PixPipeline, DDPMScheduler
28 from transformers import CLIPProcessor, CLIPModel
29 from torch.optim import AdamW
30 from torch.optim.lr_scheduler import CosineAnnealingLR
31
32 import lpips
33 from pytorch_msssim import ssim as ssim_torch
34
35 # for replication can remove
36 SEED = 42
37 random.seed(SEED)
38 np.random.seed(SEED)
39 torch.manual_seed(SEED)
40 torch.cuda.manual_seed_all(SEED)
41
42 device = "cuda" if torch.cuda.is_available() else "cpu"
43 assert device == "cuda", "This ablation requires a CUDA-enabled runtime."
44 print(f"Using device: {device}")
45
```

Using device: cuda

```
1 clip_model = CLIPModel.from_pretrained("openai/clip-vit-large-patch14").to(device)
2 clip_processor = CLIPProcessor.from_pretrained("openai/clip-vit-large-patch14")
3
4 lpips_model = lpips.LPIPS(net='vgg').to(device)
5 lpips_model.eval()
6
7 _to_tensor = transforms.ToTensor()
8
9 def pil_to_tensor(image: Image.Image) -> torch.Tensor:
10     return _to_tensor(image).unsqueeze(0).to(device)
```

```

11
12 def _match_image_size(candidate: Image.Image, reference: Image.Image) -> Image.Image:
13     if candidate.size != reference.size:
14         return candidate.resize(reference.size, Image.Resampling.LANCZOS)
15     return candidate
16
17 def calculate_lpips(img1: Image.Image, img2: Image.Image) -> float:
18     img2_aligned = _match_image_size(img2, img1)
19     x = 2 * pil_to_tensor(img1) - 1.0
20     y = 2 * pil_to_tensor(img2_aligned) - 1.0
21     with torch.no_grad():
22         return lpips_model(x, y).item()
23
24 def calculate_ssim(img1: Image.Image, img2: Image.Image) -> float:
25     img2_aligned = _match_image_size(img2, img1)
26     x = pil_to_tensor(img1)
27     y = pil_to_tensor(img2_aligned)
28     return ssim_torch(x, y, data_range=1.0, size_average=True).item()
29
30 def calculate_clip_score(image: Image.Image, text: str) -> float:
31     inputs = clip_processor(text=[text], images=image, return_tensors="pt", padding=True).to(device)
32     with torch.no_grad():
33         outputs = clip_model(**inputs)
34         image_embeds = outputs.image_embeds / outputs.image_embeds.norm(dim=-1, keepdim=True)
35         text_embeds = outputs.text_embeds / outputs.text_embeds.norm(dim=-1, keepdim=True)
36         clip_score = (image_embeds @ text_embeds.T).item()
37     return clip_score * 100
38
39 def calculate_clip_directional_similarity(original_img, edited_img, text):
40     with torch.no_grad():
41         inputs_orig = clip_processor(images=original_img, return_tensors="pt").to(device)
42         orig_embeds = clip_model.get_image_features(**inputs_orig)
43         orig_embeds = orig_embeds / orig_embeds.norm(dim=-1, keepdim=True)
44
45         inputs_edit = clip_processor(images=edited_img, return_tensors="pt").to(device)
46         edit_embeds = clip_model.get_image_features(**inputs_edit)
47         edit_embeds = edit_embeds / edit_embeds.norm(dim=-1, keepdim=True)
48
49         inputs_text = clip_processor(text=[text], return_tensors="pt", padding=True).to(device)
50         text_embeds = clip_model.get_text_features(**inputs_text)
51         text_embeds = text_embeds / text_embeds.norm(dim=-1, keepdim=True)
52
53         img_direction = edit_embeds - orig_embeds
54         img_direction = img_direction / (img_direction.norm(dim=-1, keepdim=True) + 1e-8)
55
56         null_text = clip_processor(text=[""], return_tensors="pt", padding=True).to(device)
57         null_embeds = clip_model.get_text_features(**null_text)
58         null_embeds = null_embeds / null_embeds.norm(dim=-1, keepdim=True)
59
60         text_direction = text_embeds - null_embeds
61         text_direction = text_direction / (text_direction.norm(dim=-1, keepdim=True) + 1e-8)
62
63         directional_sim = (img_direction @ text_direction.T).item()
64     return directional_sim * 100
65
66 print("metric stack ready.")
67

```

/usr/local/lib/python3.12/dist-packages/huggingface_hub/utils/_auth.py:104: UserWarning:
Error while fetching `HF_TOKEN` secret value from your vault: 'Requesting secret HF_TOKEN timed out. Secrets can only
You are not authenticated with the Hugging Face Hub in this notebook.
If the error persists, please let us know by opening an issue on GitHub (https://github.com/huggingface/huggingface_hub)
warnings.warn(
Using a slow image processor as `use_fast` is unset and a slow processor was saved with this model. `use_fast=True` w
Setting up [LPIPS] perceptual loss: trunk [vgg], v[0.1], spatial [off]
/usr/local/lib/python3.12/dist-packages/torchvision/models/_utils.py:208: UserWarning: The parameter 'pretrained' is
warnings.warn(
/usr/local/lib/python3.12/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight
warnings.warn(msg)
Loading model from: /usr/local/lib/python3.12/dist-packages/lpips/weights/v0.1/vgg.pth
Metric stack ready.

```

1 FASHIONPEDIA_STREAM_LIMIT = 100
2 MAX_TRIPLETS = 1500
3 IMAGE_SIZE = 512
4 BATCH_SIZE = 1
5 CHECKPOINT_DIR = "ablation_checkpoints"
6 os.makedirs(CHECKPOINT_DIR, exist_ok=True)
7
8 print("Loading FashionPedia stream (for qualitative eval)...")
9 dataset_fashionpedia = load_dataset("detection-datasets/fashionpedia", split="train", streaming=
10 fashionpedia_samples = list(dataset_fashionpedia.take(FASHIONPEDIA_STREAM_LIMIT))
11 print(f"Loaded {len(fashionpedia_samples)} FashionPedia samples for eval.")
12
13 print("Loading DeepFashion2 (with masks) ...")
14 deepfashion_dataset = load_dataset("SaffalPoosh/deepFashion-with-masks")
15 print(deepfashion_dataset)
16
17
18 def generate_instruction(caption1: str, caption2: str) -> str:
19     part2 = caption2.lower().split("in ")[-1] if "in " in caption2.lower() else caption2.lower()
20     return f"change to {part2}"
21
22
23 def build_triplets(hf_dataset, max_triplets: int):
24     by_pid = defaultdict(list)
25     for i in range(len(hf_dataset["train"])):
26         by_pid[hf_dataset["train"][i]["pid"]].append(i)
27
28     triplets = []
29     for pid, indices in tqdm(by_pid.items(), desc="Building triplets"):
30         if len(indices) < 2:
31             continue
32         for i in range(len(indices)):
33             for j in range(len(indices)):
34                 if i == j:
35                     continue
36                 triplets.append({"src_idx": indices[i], "tgt_idx": indices[j]})
37                 if len(triplets) >= max_triplets:
38                     return triplets
39     return triplets
40
41 triplets = build_triplets(deepfashion_dataset, max_triplets=MAX_TRIPLETS)
42 print(f"Total triplets for ablation: {len(triplets)}")
43
44
45 class TripletDataset(Dataset):
46     def __init__(self, triplet_list, hf_data, size=IMAGE_SIZE):
47         self.data = triplet_list
48         self.hf = hf_data
49         self.size = size
50         self.transform = transforms.Compose([
51             transforms.Resize((size, size), interpolation=transforms.InterpolationMode.BILINEAR),
52             transforms.ToTensor(),
53             transforms.Normalize([0.5, 0.5, 0.5], [0.5, 0.5, 0.5])
54         ])
55         self.mask_transform = transforms.Compose([
56             transforms.Resize((size, size), interpolation=transforms.InterpolationMode.NEAREST),
57             transforms.ToTensor(),
58         ])
59         self.color_jitter = transforms.ColorJitter(0.1, 0.1, 0.1, 0.05)
60
61     def __len__(self):
62         return len(self.data)
63
64     def __getitem__(self, idx):
65         t = self.data[idx]
66         src = self.hf["train"][t["src_idx"]]
67         tgt = self.hf["train"][t["tgt_idx"]]
68
69         orig_img = src["images"].convert("RGB")

```

```

70     edit_img = tgt["images"].convert("RGB")
71     src_mask = src["mask"].convert("L")
72     instruction = generate_instruction(src["caption"], tgt["caption"])
73
74     if torch.rand(1).item() < 0.5:
75         orig_img = TF.hflip(orig_img)
76         edit_img = TF.hflip(edit_img)
77         src_mask = TF.hflip(src_mask)
78
79     orig_img = self.color_jitter(orig_img)
80     edit_img = self.color_jitter(edit_img)
81
82     conditioning_image = self.transform(orig_img)
83     pixel_values = self.transform(edit_img)
84     mask_tensor = self.mask_transform(src_mask)
85     mask_tensor = (mask_tensor > 0.5).float()
86
87     return {
88         "conditioning_image": conditioning_image,
89         "pixel_values": pixel_values,
90         "prompt": instruction,
91         "mask": mask_tensor,
92     }
93
94 full_ds = TripletDataset(triplets, deepfashion_dataset, size=IMAGE_SIZE)
95 train_size = int(len(full_ds) * 0.8)
96 val_size = int(len(full_ds) * 0.1)
97 test_size = len(full_ds) - train_size - val_size
98
99 train_ds, val_ds, test_ds = torch.utils.data.random_split(
100     full_ds,
101     [train_size, val_size, test_size],
102     generator=torch.Generator().manual_seed(SEED)
103 )
104
105 train_loader = DataLoader(train_ds, batch_size=BATCH_SIZE, shuffle=True, num_workers=0)
106 val_loader = DataLoader(val_ds, batch_size=BATCH_SIZE, shuffle=False, num_workers=0)
107 test_loader = DataLoader(test_ds, batch_size=BATCH_SIZE, shuffle=False, num_workers=0)
108
109 print(f"Train: {len(train_ds)}, Val: {len(val_ds)}, Test: {len(test_ds)}")
110
111 peek = next(iter(train_loader))
112 print(f"Conditioning shape: {peek['conditioning_image'].shape}")
113 print(f"Sample prompt: {peek['prompt'][0]}")
114

```

```

Loading FashionPedia stream (for qualitative eval)...
Loaded 100 FashionPedia samples for eval.
Loading DeepFashion2 (with masks) ...
data/train-00002-of-00004-f11f7b29a225c1(...): 0%|          | 0.00/361M [00:00<?, ?B/s]
data/train-00003-of-00004-e6de6e7287b28c(...): 0%|          | 0.00/361M [00:00<?, ?B/s]
Generating train split: 0%|          | 0/40658 [00:00<?, ? examples/s]
DatasetDict({
  train: Dataset({
    features: ['images', 'gender', 'pose', 'cloth_type', 'pid', 'caption', 'mask', 'mask_overlay'],
    num_rows: 40658
  })
})
Building triplets: 0%|          | 0/7557 [00:00<?, ?it/s]
Total triplets for ablation: 1500
Train: 1200, Val: 150, Test: 150
Conditioning shape: torch.Size([1, 3, 512, 512])
Sample prompt: change to a white shirt and jeans

```

```

1 ABLATION_EPOCHS = 2
2 EVAL_PROMPTS = [
3     "change to a floral pattern",
4     "make the shirt red",
5     "add stripes to the clothes",
6     "change shirt to blue",
7 ]
8 NUM_EVAL_IMAGES = 2

```

```

9 FASHIONPEDIA_EVAL_IMAGES = 2
10
11 @dataclass
12 class AblationConfig:
13     name: str
14     learning_rate: float
15     lambda_preserve: float
16     guidance_scale: float
17     image_guidance_scale: float
18     num_inference_steps: int
19     unfreeze_strategy: str
20     weight_decay: float = 0.01
21     max_grad_norm: float = 1.0
22
23
24 def config_slug(name: str) -> str:
25     base = name.lower().replace(" ", "_")
26     return re.sub(r"^[a-z0-9_]+", "", base) or "config"
27
28
29 def apply_unfreeze_strategy(unet, strategy: str):
30     for param in unet.parameters():
31         param.requires_grad = False
32
33     trainable_names = []
34
35     def is_attn_proj(name):
36         targets = ["to_q", "to_k", "to_v", "to_out"]
37         return any(t in name for t in targets)
38
39     for name, param in unet.named_parameters():
40         should_train = False
41         if "attn" in name and is_attn_proj(name):
42             should_train = True
43         if "mid_block" in name and strategy in {"attn_mid", "attn_mid_up", "attn_mid_resnet"}:
44             should_train = True
45         if strategy in {"attn_mid_up", "attn_mid_resnet"} and "up_blocks" in name:
46             if "resnets" in name and "conv" in name:
47                 should_train = True
48         if strategy == "attn_mid_resnet" and "down_blocks" in name and "resnets" in name and "conv" in name:
49             should_train = True
50
51     if should_train:
52         param.requires_grad = True
53         trainable_names.append(name)
54
55     return trainable_names
56
57
58 def prepare_components(config: AblationConfig):
59     pipe = StableDiffusionInstructPix2PixPipeline.from_pretrained(
60         "timbrooks/instruct-pix2pix",
61         torch_dtype=torch.float16,
62         safety_checker=None,
63     ).to(device)
64     pipe.scheduler = DDPMSScheduler.from_config(pipe.scheduler.config)
65
66     unet = pipe.unet
67     tokenizer = pipe.tokenizer
68     text_encoder = pipe.text_encoder
69     vae = pipe.vae
70     noise_scheduler = pipe.scheduler
71
72     trainable_names = apply_unfreeze_strategy(unet, config.unfreeze_strategy)
73     trainable_params = [p for p in unet.parameters() if p.requires_grad]
74     optimizer = AdamW(trainable_params, lr=config.learning_rate, eps=1e-4, weight_decay=config.weight_decay)
75     scheduler = CosineAnnealingLR(optimizer, T_max=ABLATION_EPOCHS * len(train_loader))
76
77     return {
78         "pipe": pipe,

```

```

79         "unet": unet,
80         "tokenizer": tokenizer,
81         "text_encoder": text_encoder,
82         "vae": vae,
83         "noise_scheduler": noise_scheduler,
84         "optimizer": optimizer,
85         "scheduler": scheduler,
86         "trainable_params": trainable_params,
87         "trainable_names": trainable_names,
88     }
89
90
91 def encode_prompt(tokenizer, text_encoder, texts):
92     tokens = tokenizer(
93         list(texts),
94         padding="max_length",
95         max_length=tokenizer.model_max_length,
96         truncation=True,
97         return_tensors="pt",
98     ).input_ids.to(device)
99     with torch.no_grad():
100         text_embeds = text_encoder(tokens)[0]
101     return text_embeds
102
103
104 def encode_vae_images(vae, images):
105     images = images.to(device, dtype=torch.float16)
106     with torch.no_grad():
107         latents = vae.encode(images).latent_dist.sample()
108     return latents * vae.config.scaling_factor
109
110
111 def train_one_epoch(config, components):
112     unet = components["unet"]
113     vae = components["vae"]
114     tokenizer = components["tokenizer"]
115     text_encoder = components["text_encoder"]
116     noise_scheduler = components["noise_scheduler"]
117     optimizer = components["optimizer"]
118     scheduler = components["scheduler"]
119     trainable_params = components["trainable_params"]
120
121     unet.train()
122     epoch_losses = []
123     progress = tqdm(train_loader, desc=f"{config.name} - train", leave=False)
124     for batch in progress:
125         pixel_values = batch["pixel_values"].to(device, dtype=torch.float16)
126         cond_images = batch["conditioning_image"].to(device, dtype=torch.float16)
127         masks = batch["mask"].to(device, dtype=torch.float16)
128         prompts = batch["prompt"]
129
130         with torch.no_grad():
131             target_latents = encode_vae_images(vae, pixel_values)
132             cond_latents = encode_vae_images(vae, cond_images)
133             mask_latent = F.interpolate(masks, size=target_latents.shape[-2:], mode="nearest")
134             mask_latent = mask_latent.expand(-1, 4, -1, -1)
135
136             noise = torch.randn_like(target_latents)
137             timesteps = torch.randint(0, noise_scheduler.config.num_train_timesteps,
138                                     (target_latents.shape[0],), device=device, dtype=torch.long)
139             noisy_latents = noise_scheduler.add_noise(target_latents, noise, timesteps)
140             model_input = torch.cat([noisy_latents, cond_latents], dim=1)
141             text_embeds = encode_prompt(tokenizer, text_encoder, prompts)
142
143             optimizer.zero_grad()
144             with torch.autocast(device_type="cuda", dtype=torch.float16):
145                 noise_pred = unet(model_input, timesteps, encoder_hidden_states=text_embeds, return_c
146
147             M_cloth = mask_latent
148             M_bg = 1.0 - mask_latent

```

```

149     L_denoise = F.mse_loss((noise_pred * M_cloth).float(), (noise * M_cloth).float())
150     L_preserve = F.mse_loss((noise_pred * M_bg).float(), (noise * M_bg).float())
151     loss = L_denoise + config.lambda_preserve * L_preserve
152
153     if not torch.isfinite(loss):
154         continue
155
156     loss.backward()
157     torch.nn.utils.clip_grad_norm_(trainable_params, config.max_grad_norm)
158     optimizer.step()
159     scheduler.step()
160
161     epoch_losses.append(loss.item())
162     progress.set_postfix({"loss": f"{loss.item():.4f}", "lr": f"{scheduler.get_last_lr()[0]:
163
164 avg_loss = float(np.mean(epoch_losses)) if epoch_losses else float("nan")
165 return avg_loss
166
167
168 def validate_one_epoch(config, components, max_batches=50):
169     unet = components["unet"]
170     vae = components["vae"]
171     tokenizer = components["tokenizer"]
172     text_encoder = components["text_encoder"]
173     noise_scheduler = components["noise_scheduler"]
174
175     unet.eval()
176     val_losses = []
177     with torch.no_grad():
178         for batch_idx, batch in enumerate(val_loader):
179             if batch_idx >= max_batches:
180                 break
181             pixel_values = batch["pixel_values"].to(device, dtype=torch.float16)
182             cond_images = batch["conditioning_image"].to(device, dtype=torch.float16)
183             masks = batch["mask"].to(device, dtype=torch.float16)
184             prompts = batch["prompt"]
185
186             target_latents = encode_vae_images(vae, pixel_values)
187             cond_latents = encode_vae_images(vae, cond_images)
188             mask_latent = F.interpolate(masks, size=target_latents.shape[-2:], mode="nearest")
189             mask_latent = mask_latent.expand(-1, 4, -1, -1)
190
191             noise = torch.randn_like(target_latents)
192             timesteps = torch.randint(0, noise_scheduler.config.num_train_timesteps,
193                                     (target_latents.shape[0],), device=device, dtype=torch.long)
194             noisy_latents = noise_scheduler.add_noise(target_latents, noise, timesteps)
195             model_input = torch.cat([noisy_latents, cond_latents], dim=1)
196             text_embeds = encode_prompt(tokenizer, text_encoder, prompts)
197
198             with torch.autocast(device_type="cuda", dtype=torch.float16):
199                 noise_pred = unet(model_input, timesteps, encoder_hidden_states=text_embeds, retu
200
201             M_cloth = mask_latent
202             M_bg = 1.0 - mask_latent
203             L_denoise = F.mse_loss((noise_pred * M_cloth).float(), (noise * M_cloth).float())
204             L_preserve = F.mse_loss((noise_pred * M_bg).float(), (noise * M_bg).float())
205             loss = L_denoise + config.lambda_preserve * L_preserve
206             val_losses.append(loss.item())
207
208     avg_val = float(np.mean(val_losses)) if val_losses else float("nan")
209     unet.train()
210     return avg_val
211
212
213 def evaluate_config(pipe, config: AblationConfig):
214     pipe.unet.eval()
215     deepfashion_rows = []
216     fashionpedia_rows = []
217     deepfashion_visuals = []
218     fashionpedia_visuals = []

```

```

219
220 for row in range(NUM_EVAL_IMAGES):
221     actual_idx = test_ds.indices[row % len(test_ds)]
222     triplet = triplets[actual_idx]
223     src = deepfashion_dataset["train"][triplet["src_idx"]]
224     original_image = src["images"]
225     if max(original_image.size) > 512:
226         original_image = original_image.copy()
227         original_image.thumbnail((512, 512), Image.Resampling.LANCZOS)
228
229     row_visual = {"original": original_image.copy(), "edits": []}
230     for prompt in EVAL_PROMPTS:
231         with torch.autocast("cuda"):
232             edited = pipe(
233                 prompt=prompt,
234                 image=original_image,
235                 num_inference_steps=config.num_inference_steps,
236                 guidance_scale=config.guidance_scale,
237                 image_guidance_scale=config.image_guidance_scale,
238             ).images[0]
239             clip_score = calculate_clip_score(edited, prompt)
240             dir_sim = calculate_clip_directional_similarity(original_image, edited, prompt)
241             lpips_val = calculate_lpips(original_image, edited)
242             ssim_val = calculate_ssim(original_image, edited)
243             deepfashion_rows.append({
244                 "config": config.name,
245                 "dataset": "DeepFashion2",
246                 "row": row,
247                 "prompt": prompt,
248                 "clip_score": clip_score,
249                 "directional_similarity": dir_sim,
250                 "lpips": lpips_val,
251                 "ssim": ssim_val,
252             })
253             row_visual["edits"].append({"prompt": prompt, "image": edited.copy(), "clip": clip_score})
254     deepfashion_visuals.append(row_visual)
255
256 for row in range(FASHIONPEDIA_EVAL_IMAGES):
257     sample = fashionpedia_samples[row]
258     original_image = sample['image']
259     if max(original_image.size) > 512:
260         original_image = original_image.copy()
261         original_image.thumbnail((512, 512), Image.Resampling.LANCZOS)
262
263     row_visual = {"original": original_image.copy(), "edits": []}
264     for prompt in EVAL_PROMPTS:
265         with torch.autocast("cuda"):
266             edited = pipe(
267                 prompt=prompt,
268                 image=original_image,
269                 num_inference_steps=config.num_inference_steps,
270                 guidance_scale=config.guidance_scale,
271                 image_guidance_scale=config.image_guidance_scale,
272             ).images[0]
273             clip_score = calculate_clip_score(edited, prompt)
274             dir_sim = calculate_clip_directional_similarity(original_image, edited, prompt)
275             lpips_val = calculate_lpips(original_image, edited)
276             ssim_val = calculate_ssim(original_image, edited)
277             fashionpedia_rows.append({
278                 "config": config.name,
279                 "dataset": "FashionPedia",
280                 "row": row,
281                 "prompt": prompt,
282                 "clip_score": clip_score,
283                 "directional_similarity": dir_sim,
284                 "lpips": lpips_val,
285                 "ssim": ssim_val,
286             })
287             row_visual["edits"].append({"prompt": prompt, "image": edited.copy(), "clip": clip_score})
288     fashionpedia_visuals.append(row_visual)

```



```

289
290     return {
291         "deepfashion_metrics": pd.DataFrame(deepfashion_rows),
292         "fashionpedia_metrics": pd.DataFrame(fashionpedia_rows),
293         "deepfashion_visuals": deepfashion_visuals,
294         "fashionpedia_visuals": fashionpedia_visuals,
295     }
296
297

```

```

1 ABLATION_CONFIGS = [
2     AblationConfig(
3         name="A:  $\lambda 0.5$ , lr1e-5, attn+mid+up",
4         learning_rate=1e-5,
5         lambda_preserve=0.5,
6         guidance_scale=7.5,
7         image_guidance_scale=1.5,
8         num_inference_steps=40,
9         unfreeze_strategy="attn_mid_up",
10    ),
11    AblationConfig(
12        name="B:  $\lambda 0.8$ , lr2e-5, attn+mid+up",
13        learning_rate=2e-5,
14        lambda_preserve=0.8,
15        guidance_scale=8.0,
16        image_guidance_scale=1.7,
17        num_inference_steps=50,
18        unfreeze_strategy="attn_mid_up",
19    ),
20    AblationConfig(
21        name="C:  $\lambda 0.3$ , lr5e-6, attn+mid",
22        learning_rate=5e-6,
23        lambda_preserve=0.3,
24        guidance_scale=6.5,
25        image_guidance_scale=1.3,
26        num_inference_steps=35,
27        unfreeze_strategy="attn_mid",
28    ),
29    AblationConfig(
30        name="D:  $\lambda 0.5$ , lr1e-5, attn+mid+res",
31        learning_rate=1e-5,
32        lambda_preserve=0.5,
33        guidance_scale=7.0,
34        image_guidance_scale=1.4,
35        num_inference_steps=40,
36        unfreeze_strategy="attn_mid_resnet",
37    ),
38 ]
39
40 ablation_runs: List[Dict[str, Any]] = []
41
42 for cfg in ABLATION_CONFIGS:
43     print(f"\n===== Running ablation: {cfg.name} =====")
44     components = prepare_components(cfg)
45     trainable_count = sum(p.numel() for p in components["trainable_params"])
46
47     history = []
48     for epoch in range(ABLATION_EPOCHS):
49         train_loss = train_one_epoch(cfg, components)
50         val_loss = validate_one_epoch(cfg, components, max_batches=len(val_loader))
51         history.append({"epoch": epoch + 1, "train_loss": train_loss, "val_loss": val_loss})
52         print(f"Epoch {epoch+1}/{ABLATION_EPOCHS} - train: {train_loss:.4f}, val: {val_loss:.4f}")
53
54         ckpt_path = os.path.join(CHECKPOINT_DIR, f"{config_slug(cfg.name)}_epoch{epoch+1}.pt")
55         torch.save({
56             "config": asdict(cfg),
57             "epoch": epoch + 1,
58             "unet_state_dict": components["unet"].state_dict(),
59             "optimizer_state_dict": components["optimizer"].state_dict(),

```

```
60         "scheduler_state_dict": components["scheduler"].state_dict(),
61     }, ckpt_path)
62     print(f" ↳ checkpoint saved to {ckpt_path}")
63
64     eval_artifacts = evaluate_config(components["pipe"], cfg)
65
66     run_summary = {
67         "config": asdict(cfg),
68         "history": history,
69         "trainable_params": trainable_count,
70         "deepfashion_metrics": eval_artifacts["deepfashion_metrics"],
71         "fashionpedia_metrics": eval_artifacts["fashionpedia_metrics"],
72         "deepfashion_visuals": eval_artifacts["deepfashion_visuals"],
73         "fashionpedia_visuals": eval_artifacts["fashionpedia_visuals"],
74     }
75     ablation_runs.append(run_summary)
76
77     del components["pipe"]
78     del components
79     torch.cuda.empty_cache()
80     gc.collect()
81
```



```

===== Running ablation: A:  $\lambda$ 0.5, lr1e-5, attn+mid+up =====
model_index.json: 0%|          | 0.00/616 [00:00<?, ?B/s]
Fetching 13 files: 0%|          | 0/13 [00:00<?, ?it/s]
config.json: 0%|          | 0.00/617 [00:00<?, ?B/s]
merges.txt: 0.00B [00:00, ?B/s]
scheduler_config.json: 0%|          | 0.00/569 [00:00<?, ?B/s]
preprocessor_config.json: 0%|          | 0.00/518 [00:00<?, ?B/s]
special_tokens_map.json: 0%|          | 0.00/472 [00:00<?, ?B/s]
tokenizer_config.json: 0%|          | 0.00/806 [00:00<?, ?B/s]
vocab.json: 0.00B [00:00, ?B/s]
text_encoder/model.safetensors: 0%|          | 0.00/492M [00:00<?, ?B/s]
config.json: 0.00B [00:00, ?B/s]
config.json: 0%|          | 0.00/553 [00:00<?, ?B/s]
vae/diffusion_pytorch_model.safetensors: 0%|          | 0.00/335M [00:00<?, ?B/s]
UNET/diffusion_pytorch_model.safetensors: 0%|          | 0.00/3.44G [00:00<?, ?B/s]
Loading pipeline components...: 0%|          | 0/6 [00:00<?, ?it/s]
`torch_dtype` is deprecated! Use `dtype` instead!
A:  $\lambda$ 0.5, lr1e-5, attn+mid+up - train: 0%|          | 0/1200 [00:00<?, ?it/s]
Epoch 1/2 - train: 0.0405, val: 0.0433
  ↳ checkpoint saved to ablation_checkpoints/a_05_lr1e5_attnmidup_epoch1.pt
A:  $\lambda$ 0.5, lr1e-5, attn+mid+up - train: 0%|          | 0/1200 [00:00<?, ?it/s]
Epoch 2/2 - train: 0.0412, val: 0.0400
  ↳ checkpoint saved to ablation_checkpoints/a_05_lr1e5_attnmidup_epoch2.pt
0%|          | 0/40 [00:00<?, ?it/s]
0%|          | 0/40 [00:00<?, ?it/s]
0%|          | 0/40 [00:00<?, ?it/s]
0%|          | 0/40 [00:00<?, ?it/s]
0%|          | 0/40 [00:00<?, ?it/s]
0%|          | 0/40 [00:00<?, ?it/s]
0%|          | 0/40 [00:00<?, ?it/s]
0%|          | 0/40 [00:00<?, ?it/s]
0%|          | 0/40 [00:00<?, ?it/s]
0%|          | 0/40 [00:00<?, ?it/s]

```

```

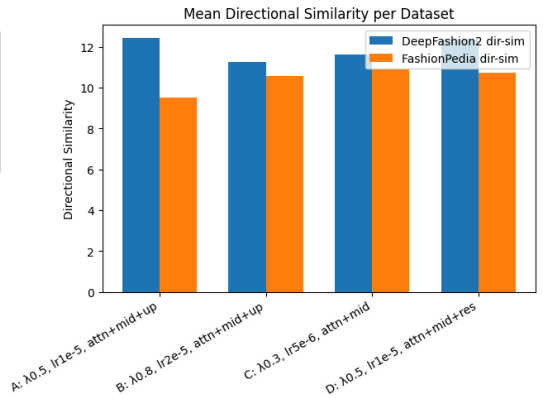
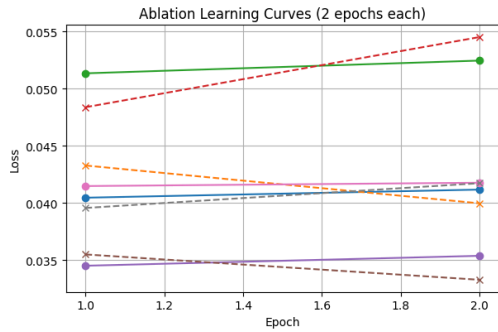
1 history_rows = []
2 all_metrics_rows = []
3 for run in ablation_runs:
4     cfg_name = run["config"]["name"]
5     for entry in run["history"]:
6         history_rows.append({"config": cfg_name, **entry})
7     all_metrics_rows.append(run["deepfashion_metrics"])
8     all_metrics_rows.append(run["fashionpedia_metrics"])
9
10 history_df = pd.DataFrame(history_rows)
11 metrics_df = pd.concat(all_metrics_rows, ignore_index=True)
12 metric_cols = ["clip_score", "directional_similarity", "lpips", "ssim"]
13 numeric_metrics = metrics_df[["config", "dataset"] + metric_cols].copy()
14
15 fig, axes = plt.subplots(1, 2, figsize=(16, 5))
16 for cfg_name, group in history_df.groupby("config"):
17     axes[0].plot(group["epoch"], group["train_loss"], marker='o', label=f"{cfg_name} - train")
18     axes[0].plot(group["epoch"], group["val_loss"], marker='x', linestyle='--', label=f"{cfg_name} - val")
19 axes[0].set_title("Ablation Learning Curves (2 epochs each)")
20 axes[0].set_xlabel("Epoch")
21 axes[0].set_ylabel("Loss")
22 axes[0].grid(True)
23 axes[0].legend(bbox_to_anchor=(1.02, 1), loc='upper left')
24
25 agg = numeric_metrics.groupby(["config", "dataset"], as_index=False)[metric_cols].mean()
26 configs = [run["config"]["name"] for run in ablation_runs]
27 x = np.arange(len(configs))
28 width = 0.35
29
30 for idx, dataset in enumerate(["DeepFashion2", "FashionPedia"]):
31     subset = agg[agg["dataset"] == dataset].set_index("config")
32     scores = [subset.loc[cfg, "directional_similarity"] for cfg in configs]
33     axes[1].bar(x + (idx - 0.5) * width, scores, width=width, label=f"{dataset} dir-sim")
34
35 axes[1].set_title("mean Directional Similarity per Dataset")
36 axes[1].set_xticks(x)
37 axes[1].set_xticklabels(configs, rotation=30, ha='right')
38 axes[1].set_ylabel("Directional Similarity")
39 axes[1].legend()
40 plt.tight_layout()
41 plt.show()
42

```

```

43 print("Metric summary (CLIP / DirSim / LPIPS / SSIM):")
44 display(agg[["config", "dataset"] + metric_cols])
45

```



Metric summary (CLIP / DirSim / LPIPS / SSIM):

	config	dataset	clip_score	directional_similarity	lpips	ssim
0	A: λ0.5, lr1e-5, attn+mid+up	DeepFashion2	20.509795	12.442044	0.355809	0.495023
1	A: λ0.5, lr1e-5, attn+mid+up	FashionPedia	19.427384	9.522973	0.408848	0.487818
2	B: λ0.8, lr2e-5, attn+mid+up	DeepFashion2	20.870804	11.254064	0.270278	0.554784
3	B: λ0.8, lr2e-5, attn+mid+up	FashionPedia	19.699326	10.573883	0.364491	0.528399
4	C: λ0.3, lr5e-6, attn+mid	DeepFashion2	21.855139	11.615444	0.331672	0.483953
5	C: λ0.3, lr5e-6, attn+mid	FashionPedia	20.084592	10.888882	0.399855	0.485995
6	D: λ0.5, lr1e-5, attn+mid+res	DeepFashion2	21.722468	12.338617	0.296990	0.492303
7	D: λ0.5, lr1e-5, attn+mid+res	FashionPedia	20.142740	10.744428	0.377625	0.515408

```

1 # == Visualization Set 2: qualitative grids for both datasets ==
2 def plot_dataset_grid(runs, key, title):
3     rows = len(runs)
4     cols = len(EVAL_PROMPTS) + 1
5     fig, axes = plt.subplots(rows, cols, figsize=(3.8 * cols, 3.2 * rows))
6     if rows == 1:
7         axes = np.expand_dims(axes, axis=0)
8
9     for r, run in enumerate(runs):
10         cfg_name = run["config"]["name"]
11         sample_row = run[key][0]
12         axes[r, 0].imshow(sample_row["original"])
13         axes[r, 0].set_title(f"{cfg_name}\noriginal", fontsize=10)
14         axes[r, 0].axis('off')
15         for c, edit in enumerate(sample_row["edits"], start=1):
16             axes[r, c].imshow(edit["image"])
17             axes[r, c].set_title(f"{edit['prompt']}\nCLIP {edit['clip']:.1f} / Dir {edit['dir']:.1f}")
18             axes[r, c].axis('off')
19     plt.suptitle(title, fontsize=14)
20     plt.tight_layout()
21     plt.show()
22
23 plot_dataset_grid(ablation_runs, "deepfashion_visuals", "Ablation Outputs – DeepFashion2 test set")
24 plot_dataset_grid(ablation_runs, "fashionpedia_visuals", "Ablation Outputs – FashionPedia validation set")

```