

# PROMO: Promptable Outfitting for Efficient High-Fidelity Virtual Try-On

## Supplementary Material

### A. Supplementary Quantitative Results

The full results of ablating **Parsing Area Loss**, **Style Prompt** and **3D-RoPE**

	DressCode						VITON-HD					
	Paired				Unpaired		Paired				Unpaired	
	SSIM ↑	LPIPS ↓	FID ↓	KID ↓	FID ↓	KID ↓	SSIM ↑	LPIPS ↓	FID ↓	KID ↓	FID ↓	KID ↓
- w/o parsing area loss	<u>0.8896</u>	<b>0.0872</b>	<b>3.2847</b>	<u>0.5076</u>	<b>4.6390</b>	0.9535	0.8401	<u>0.1342</u>	<u>7.0419</u>	<b>1.2711</b>	<b>9.4283</b>	<b>1.2711</b>
- w/o style prompt	<u>0.8896</u>	0.0926	3.7178	0.8926	<u>5.3516</u>	<u>0.6180</u>	<u>0.8432</u>	0.1376	7.8145	1.5427	10.1242	<u>1.5575</u>
- w/o 3D-RoPE	0.8702	0.1304	6.7293	1.7160	7.8175	2.2801	0.8335	0.1609	10.1979	1.5219	11.8961	2.2801
PROMO (Ours)	<b>0.8913</b>	<u>0.0887</u>	<u>3.3103</u>	<b>0.4902</b>	<u>5.3516</u>	<b>0.4992</b>	<b>0.8619</b>	<b>0.1111</b>	<b>6.8903</b>	<u>1.4895</u>	<u>9.9034</u>	1.9174

Table 1. Ablation study on DressCode and VITON-HD datasets. The best and second-best results are highlighted in **bold** and underlined, respectively.

### B. Implementation Details of TryOff Model

We trained our TryOff model based on FLUX.1 Kontext [dev] [3], using LoRA [2] with rank of 16 and alpha of 4. We employed the AdamW optimizer [4] with a learning rate of 1e-4. The training dataset consists of person-garment pairs from DressCode [5] and Viton-HD [1] datasets.

For garment extraction, we defined category-specific prompts as follows:

```
PROMPT_DICT = {
    "lower_body": "extract only the lower body garment over a white
                  background, product photography style",
    "upper_body": "extract only the upper body garment over a white
                  background, product photography style",
    "dresses": "extract only the dresses garment over a white
               background, product photography style",
}
```

### C. Prompt Extraction Details

The following code is the complete pydantic schema, and we use the OpenAPI json formatted schema description generated by pydantic using `Output.model_json_schema()` to guide the LLM to output the expected result.

```
1 from pydantic import BaseModel
2
3 class Person(BaseModel):
4     body_shape: str = Field(description="Body type in 1 word", examples=["slim",
5     ↪ "athletic", "curvy", "plus-size", "average"])
6     gender: Literal["man", "woman"]
7     hair_length: str = Field(description="Hair length in 2-3 words", examples=["short",
8     ↪ "shoulder-length", "long", "very-long"])
9     pose: str = Field(description="Overall body pose in 5-10 words")
10    hand_pose: str = Field(description="Hand/arm position in 5-10 words")
11
12 class GarmentBase(BaseModel):
```

```

11 garment_type: str = Field(description="Describe specific garment type in 1 words.",
12     ↪ examples=["t-shirt", "sweater", "jeans", "skirt", "shorts", "dress", "romper"])
13 color: str = Field(description="Describe cloth color in 1 words.")
14 material: str = Field(description="Describe material in 1 words.")
15 fit: str = Field(description="Describe fit of the garment in 1 words.",
16     ↪ examples=["tight", "loose", "relaxed", "fitted", "oversized"])
17
18 class GarmentUpperBase(GarmentBase):
19     neckline: str = Field(description="Describe neckline in 1 word.")
20     sleeve_length: Literal["sleeveless", "short sleeve (above elbow)", "elbow length",
21     ↪ "3/4 sleeve (mid forearm)", "long sleeve (at wrist)", "extra long (below wrist,
22     ↪ covering hand)"]
23     sleeve_rolling_style: Literal["sleeveless", "rolled down", "rolled up", "rolled in
24     ↪ the middle"]
25
26 class GarmentUpper(GarmentUpperBase):
27     type: Literal["upper_body"]
28     upper_garment_length: Literal["cropped length", "waist level", "below waist, covering
29     ↪ lower garment"]
30
31 class GarmentLowerBase(GarmentBase):
32     lower_garment_length: Literal["out of frame", "short (above knee)", "knee-length
33     ↪ (precisely reaches knee)", "midi (mid-calf)", "ankle-length", "floor-length"]
34
35 class GarmentLower(GarmentLowerBase):
36     type: Literal["lower_body"]
37
38 class GarmentFull(GarmentUpperBase, GarmentLowerBase):
39     type: Literal["full_body"]
40
41 class GarmentShoes(GarmentBase):
42     type: Literal["shoes"]
43     shoes_length: Literal["above knee", "knee height", "mid-calf", "ankle height",
44     ↪ "low-top", "no shoes"]
45
46 class GarmentSocks(GarmentBase):
47     type: Literal["socks"]
48     socks_length: Literal["above knee", "knee height", "mid-calf", "ankle height",
49     ↪ "low-top", "no socks"]
50
51 class Output(BaseModel):
52     person: Person
53     garments: list[Annotated[GarmentUpper | GarmentLower | GarmentFull | GarmentShoes |
54     ↪ GarmentSocks, Field(discriminator="type")]]

```

## References

- [1] Seunghwan Choi, Sunghyun Park, Minsoo Lee, and Jaegul Choo. Viton-hd: High-resolution virtual try-on via misalignment-aware normalization, 2021. [1](#)
- [2] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. [1](#)
- [3] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, Sumith Kulal, Kyle Lacey, Yam Levi, Cheng Li, Dominik Lorenz, Jonas Müller, Dustin Podell, Robin Rombach, Harry Saini, Axel Sauer, and Luke Smith. Flux.1 kontext: Flow matching for in-context image generation and editing in latent space, 2025. [1](#)
- [4] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. [1](#)
- [5] Davide Morelli, Matteo Fincato, Marcella Cornia, Federico Landi, Fabio Cesari, and Rita Cucchiara. Dress code: High-resolution multi-category virtual try-on, 2022. [1](#)