

Cinderella: Advanced Virtual Try-on

Anonymous CVPR submission

Paper ID *****

Abstract

001 *E-commerce platforms can generate virtual try-on images*
002 *but cannot accurately represent garment fit variations (tight*
003 *vs. loose) at the same detail fidelity required for purchase*
004 *decisions. Our approach combines IDM-VTON’s dual-*
005 *encoding architecture (IP-Adapter for semantics, Garment-*
006 *Net for low-level details) with a novel size-conditioning*
007 *module that learns garment-to-body dimensional ratios,*
008 *producing per-image try-on results where an XL garment*
009 *appears appropriately loose on a small frame. Fine-tuning*
010 *a fashion-specific IP-Adapter and integrating spatial size*
011 *guidance maps addresses the core challenge: preserving in-*
012 *ntricate garment patterns while respecting physical size con-*
013 *straints. The deliverable includes a 2D try-on pipeline with*
014 *optional multi-view 3D rendering, comprehensive evalua-*
015 *tion metrics, and an interactive demo for real-world valida-*
016 *tion. We propose a comprehensive virtual try-on addressing*
017 *three aspects: detail preservation, size aware generation*
018 *and 3D controllable rendering to address the critical limi-*
019 *tations*

1. Introduction

020 When shoppers evaluate garments online, they face a critical question that existing virtual try-on (VTON) systems
021 fail to answer: *will this size actually fit me?* Current VTON
022 methods excel at warping garments to body contours but ig-
023 nore the fundamental relationship between garment dimen-
024 sions and body proportions—an XL shirt renders identically
025 fitted whether the model is petite or plus-size. This limita-
026 tion stems from treating try-on purely as a texture trans-
027 fer problem rather than a size-aware transformation task.
028 Our goal is a learning-based system that generates authentic
029 virtual try-on images with explicit size control, where gar-
030 ment dimensions determine visual fit characteristics (drap-
031 ing, tightness, coverage area). Each output explains how the
032 size was determined through interpretable size maps and at-
033 tention visualizations, enabling both automated sizing rec-
034 ommendations and human verification.
035

2. Limits of Current Practice

037

Today’s VTON systems face three persistent challenges. **First, size blindness:** models like HR-VITON, VITON-HD, and even recent diffusion methods (LaDI-VTON, DCI-VTON, StableVITON) warp garments to fit target poses without considering whether a garment is XS or XXL—the same hoodie appears perfectly fitted on all body types. **Second, detail loss:** while GAN-based methods struggle with pattern preservation, diffusion models often blur logos, text, and intricate designs when conditioning is insufficient. **Third, lack of controllability:** users cannot specify “show me this in a size larger” without re-running with different garment images. Prior work addresses detail (IDM-VTON’s GarmentNet) or size (COTTON’s landmark-based warping) separately, but no system combines both with diffusion model quality. This gap motivates our integrated design.

038

039

040

041

042

043

044

045

046

047

048

049

050

051

052

053

3. Approach and Key Insight

054

We treat size-aware try-on as a multi-modal conditioning problem: garment appearance (preserved via dual encoding), body geometry (captured via DensePose and pose key-points), and size relationship (learned from garment-body dimension ratios). By injecting size information into attention mechanisms rather than relying solely on implicit warping, the model learns when to generate tight-fitting versus loose-draping results. Fine-tuning IP-Adapter on fashion-specific semantics (materials, styles, necklines) further ensures that “cotton t-shirt” differs meaningfully from “silk blouse” in rendered texture.

055

056

057

058

059

060

061

062

063

064

065

3.1. Data and Size Annotation Policy

066

Base Datasets. VITON-HD (11,647 pairs, 1024×768) for primary training; DressCode (multi-category) for cross-category generalization testing; DeepFashion2 (801K images, 13 categories, 491K clothing items with bounding boxes, dense landmarks, and commercial consumer-to-shop pairs) for fashion-specific semantic understanding and landmark detection pre-training; custom In-the-Wild set (1,500+ pairs) with size annotations for size-aware eval-

067

068

069

070

071

072

073

074

075 uation. DeepFashion2’s rich annotations including cloth-
 076 ing landmarks, scale, occlusion, zoom-in, viewpoint, and
 077 bounding boxes provide ideal training data for our 10-point
 078 garment landmark predictor and material/style classifica-
 079 tion components.

080 **Size Labeling.** We extract garment dimensions (shoulder
 081 width, torso/sleeve length) via 10-point CNN land-
 082 mark predictor pre-trained on DeepFashion2 and fine-tuned
 083 on 500 VITON-HD samples, and body dimensions from
 084 pose keypoints (OpenPose/MediaPipe). Size ratio $r =$
 085 (garment/body width, garment/body length) maps to dis-
 086 crete labels: tight ($r < 0.9$), fitted ($0.9 \leq r < 1.1$), loose
 087 ($1.1 \leq r < 1.3$), oversized ($r \geq 1.3$). 500 images man-
 088 ually verified in Label Studio; augmentation scales garments
 089 0.7–1.5× to balance distribution.

090 3.2. Preprocessing

091 **Person Image Processing.** (1) Human parsing:
 092 SCHP/Graphonomy produces 20-class segmentation.
 093 (2) DensePose estimation: Detectron2 extracts UV body
 094 surface maps. (3) Pose keypoints: OpenPose provides an
 095 18-point skeleton. (4) Agnostic mask: Remove original
 096 garment while preserving arms (to retain skin tone, tattoos,
 097 width). (5) Encoding: All inputs (masked person, Dense-
 098 Pose, mask) passed through VAE encoder to latent space
 099 ($128 \times 96 \times 4$).

100 **Garment Image Processing.** (1) Segmentation: SAM/U²-
 101 Net isolates garment from background. (2) Landmark de-
 102 tection: 10-point predictor (neck, shoulders, elbows, wrists,
 103 hips) trained on a curated set. (3) Caption generation:
 104 Fashion-specific CLIP tagger extracts: sleeve type, neck-
 105 line, material, features → Template: “A photo of [at-
 106 tributes]”. (4) Encoding: Garment latent ($128 \times 96 \times 4$)
 107 + CLIP image features (257×1280).

108 3.3. Model Architecture and Training

109 **IDM-VTON Core (Base).** Our base architecture consists
 110 of three components: *TryonNet*, an SDXL Inpainting UNet
 111 with 13-channel input (noised latent + mask + masked per-
 112 son + DensePose); *GarmentNet*, a frozen SDXL UNet en-
 113 coder extracting multi-scale garment features (low-level de-
 114 tails: patterns, logos, textures); and *IP-Adapter*, combining
 115 frozen CLIP ViT-H/14 with trainable projection layers for
 116 high-level semantics (style, color, material).

117 **Size Module (Novel).** We introduce two novel components
 118 for size awareness: *Size Encoder*, an MLP mapping size ra-
 119 tio [width_ratio, length_ratio, sleeve_ratio] → 768-dim em-
 120 bedding; and *Size Controller*, a CNN generating spatial size
 121 maps (per-pixel tight/loose guidance) from fused person +
 122 garment + size features.

123 **Size-Aware Attention.** We modify self/cross-attention lay-
 124 ers to incorporate size information. *Self-attention* modu-
 125 lates token importance based on size maps (loose regions

126 attend more to garment structure), while *cross-attention* in-
 127 jects size embeddings alongside text/image features.

128 3.4. Training Protocol

129 **Stage 1: Base IDM-VTON.** We use pretrained TryonNet
 130 and fine-tune IP-Adapter projection on VITON-HD while
 131 freezing GarmentNet to preserve SDXL prior knowledge.

132 **Stage 2: IP-Adapter Fashion Fine-tuning (30 epochs).**
 133 We freeze TryonNet/GarmentNet and fine-tune IP-Adapter
 134 on 5,000 fashion images, adding custom attention pro-
 135 cessors with garment-specific gating.

136 **Stage 3: Size Module Training (50 epochs).** We freeze the
 137 base IDM-VTON and train the size encoder/controller using
 138 an augmented dataset with synthetic size variations. The
 139 loss function combines reconstruction, size consistency, and
 140 spatial alignment: $\mathcal{L}_{\text{rec}} + 0.5 \cdot \mathcal{L}_{\text{size_consistency}} + 0.3 \cdot \mathcal{L}_{\text{spatial}}$.

141 **Stage 4: Joint Fine-tuning (30 epochs).** We unfreeze the
 142 TryonNet decoder, IP-Adapter, and size modules for end-
 143 to-end optimization with multi-objective loss: $0.3 \cdot \mathcal{L}_{\text{idm}} +$
 144 $0.25 \cdot \mathcal{L}_{\text{ip}} + 0.25 \cdot \mathcal{L}_{\text{size}} + 0.15 \cdot \mathcal{L}_{\text{detail}} + 0.05 \cdot \mathcal{L}_{\text{human}}$.

145 **Stage 5: GS-VTON 3D (20 epochs, optional).** We op-
 146 tionally use 2D outputs as supervision for 3D Gaussian
 147 splatting scenes, optimizing Gaussian parameters (position,
 148 color, opacity) to match multi-view try-on results.

149 4. Implementation Details

150 **Hardware and Training Time.** Our system requires $4 \times$
 151 A100 (80GB) or H100 GPUs with total training time of ap-
 152 proximately 70 hours: Base (25h), IP-Adapter (15h), Size
 153 (20h), and Joint fine-tuning (10h).

154 **Development Timeline.** Week 1: Data curation (size anno-
 155 tations, landmark training), infrastructure setup, and base-
 156 line IDM-VTON implementation. Week 2: IP-Adapter
 157 fashion fine-tuning with custom processors and material
 158 classification. In parallel, size module development (en-
 159 coder, controller, attention modifications) and integration
 160 testing. Weeks 3–4: Joint fine-tuning and ablation studies
 161 (w/o size, w/o IP-Adapter tuning, w/o GarmentNet). Week
 162 5: Optional GS-VTON 3D extension and multi-view con-
 163 sistency. Week 6: Evaluation, user studies, and demo devel-
 164 opment (Streamlit interface).

165 5. Evaluation and Success Criteria

166 5.1. Mid-Term Evaluation

167 We evaluate our model on held-out validation scenes
 168 (DressCode + 200 In-the-Wild images) using three criteria:

169 **Detail Preservation.** We require LPIPS < 0.12 compared
 170 to the IDM-VTON baseline (≈ 0.102) to ensure our size-
 171 aware modifications do not degrade visual quality.

172 **Size Awareness Baseline.** We compare against two base-
 173 lines: (a) IDM-VTON without size conditioning, and (b) a

Table 1. Quantitative evaluation metrics and targets. Our model aims to match or exceed the baseline IDM-VTON while introducing size awareness capabilities.

Metric	Target	Baseline (IDM-VTON)
LPIPS (Detail)	< 0.10	0.102
SSIM (Structure)	> 0.90	0.870
FID (Realism)	< 6.0	6.29
Size Accuracy	> 85%	NA
CLIP-I (Similarity)	> 0.90	0.883
DISTS	> 0.85	—
LPIPS-Clothing	< 0.08	—
LPIPS-Person	< 0.12	—
mIoU	> 0.75	—
FTS	> 0.80	—
Size Accuracy (Ratio)	> 85%	—

174 size-only heuristic that scales garments by dimension ratio
 175 and performs simple paste operations.

176 **Qualitative Attention Verification.** Using Grad-CAM vi-
 177 sualization, we verify that attention mechanisms focus on
 178 garment boundaries rather than background for 90%+ of
 179 samples, as rated by two annotators on 100 patches.

180 5.2. Final Evaluation

181 **Quantitative Metrics.** We evaluate on a test set of 500
 182 images using the following metrics:

183 References

- 184 [1] Yukang Cao, Masoud Hadi, Liang Pan, and Ziwei Liu. Con-
 185 trollable 3d virtual try-on with gaussian splatting, 2024. arXiv
 186 preprint arXiv:2410.05259.
- 187 [2] Chieh-Yun Chen, Yi-Chung Chen, Hong-Han Shuai, and
 188 Wen-Huang Cheng. Size does matter: Size-aware virtual try-
 189 on via clothing-oriented transformation. In *ICCV*, 2023.
- 190 [3] Yisol Choi, Sangkyung Kwak, Kyungmin Lee, Hyungwon
 191 Choi, and Jinwoo Shin. Improving diffusion models for au-
 192 thentic virtual try-on in the wild. In *ECCV*, 2024.
- 193 [4] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao,
 194 Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer White-
 195 head, Alexander Berg, Wan-Yen Lo, Piotr Dollár, and Ross
 196 Girshick. Segment anything. In *ICCV*, 2023.
- 197 [5] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann,
 198 Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach.
 199 Sdxl: Improving latent diffusion models for high-resolution
 200 image synthesis, 2023. arXiv preprint arXiv:2307.01952.
- 201 [6] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
 202 Patrick Esser, and Björn Ommer. High-resolution image syn-
 203 thesis with latent diffusion models. In *CVPR*, 2022.
- 204 [7] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang.
 205 Ip-adapter: Text compatible image prompt adapter for
 206 text-to-image diffusion models, 2023. arXiv preprint
 207 arXiv:2308.06721.
- 208 [8] Shufang Zhang, Hang Qian, Minxue Ni, Yaxuan Li and
 209 Wenxin Ding, and Jun Liu. Size variable virtual try-on with
 210 diffusion models, 2025. arXiv preprint arXiv:2504.00562.