

AI Based Virtual Try-On System Using DL

Aditya Kawade¹, Saurabh Brahmkar², Vishal Dipake³, Sagar Kharat⁴, Prof. S. B. Nimbekar⁵

^{1, 2, 3, 4}Dept of Computer Engineering

⁵professor, Dept of Computer Engineering

Abstract- *The proliferation of e-commerce has fundamentally altered consumer purchasing behaviour, yet a persistent challenge remains: the inability to physically try on garments before purchase. This paper presents a Virtual Clothing Try-On (VCTO) system that bridges this gap by combining a custom Conditional Generative Adversarial Network (cGAN) with a Geometric Matching Module (GMM) embedded in a full-stack MERN (MongoDB, Express.js, React.js, Node.js) web application. The proposed system accepts a reference person image and a desired garment image as inputs and synthesises a photorealistic composite image by (i) estimating human body pose using a lightweight keypoint detector, (ii) warping the garment via Thin-Plate Spline (TPS) transformation, and (iii) generating the final try-on image through an adversarial training scheme. Experimental evaluation on the VITON-HD benchmark dataset yields a Structural Similarity Index (SSIM) of 0.873, Fréchet Inception Distance (FID) of 8.34, and Learned Perceptual Image Patch Similarity (LPIPS) of 0.072, outperforming several baseline GAN-based methods. The system achieves an average inference latency of 320 ms on a single NVIDIA RTX 3060 GPU, making it suitable for near-real-time web deployment.*

Keywords: Virtual Try-On, Conditional GAN, Geometric Matching Module, Thin-Plate Spline, MERN Stack, Deep Learning, Image Synthesis, Pose Estimation.

I. INTRODUCTION

Online fashion retail has experienced unprecedented growth over the past decade, with global revenues projected to exceed USD 1.2 trillion by 2027. Despite this growth, online apparel shopping continues to suffer from one of the highest return rates across all product categories, primarily because consumers cannot assess garment fit, drape, or colour fidelity prior to purchase. Studies indicate that approximately 30–40% of all fashion items purchased online are returned, imposing substantial logistical and financial costs on retailers.

Virtual try-on (VTO) technology offers a compelling solution by enabling consumers to visualise how a garment would appear on their body without physical interaction. Early approaches relied on 3D body avatar reconstruction, which required specialised hardware and time-consuming scanning processes. More recently, image-based 2D virtual try-on systems leveraging deep generative models — particularly

Generative Adversarial Networks (GANs) and, latterly, diffusion models — have emerged as practical and scalable alternatives.

This paper proposes a holistic Virtual Clothing Try-On System (VCTOS) that integrates a custom-trained cGAN pipeline with a production-grade MERN stack web platform. Our specific contributions are:

- A custom Conditional GAN model with a U-Net generator and PatchGAN discriminator tailored for garment-person image synthesis.
- A Geometric Matching Module (GMM) employing Thin-Plate Spline (TPS) transformation to warp garments accurately to body shape.
- Integration with a lightweight MediaPipe-based pose estimator for real-time keypoint extraction.
- A scalable MERN stack web application with RESTful API integration, enabling browser-based virtual try-on with sub-400 ms end-to-end latency.
- Quantitative evaluation on the VITON-HD dataset with comparison against VITON, CP-VTON, ACGPN, and HR-VITON baselines.

II. LITERATURE REVIEW

2.1 Early Image-Based Virtual Try-On

Han et al. [1] introduced VITON, the first image-based virtual try-on system, which decomposed the problem into a clothing-agnostic person representation and a coarse-to-fine synthesis pipeline. Wang et al. [2] extended this work with CP-VTON, incorporating a Geometric Matching Module trained end-to-end to improve garment warping accuracy. These foundational systems established the benchmark pipeline that subsequent works have refined.

2.2 GAN-Based Advancements

Yang et al. [3] proposed ACGPN (Adaptive Content Generating and Preserving Network), which adaptively preserves or regenerates image content according to a semantic layout. Islam et al. [4] conducted a comprehensive survey of deep learning methods for virtual try-on, categorising approaches into image-based, multi-pose, and video try-on paradigms and noting that cGAN architectures

remain the dominant approach for image-based systems due to their balance of quality and inference speed.

2.3 High-Resolution and Multi-Pose Systems

Choi et al. [5] introduced HR-VITON, achieving high-resolution try-on by disentangling misalignment and occlusion handling conditions. StyleVTON [6] extended multi-pose capability while preserving clothing texture details through a dual U-Net + pix2pix discriminator framework. More recently, diffusion-based approaches such as LaDI-VTON [7] have demonstrated superior FID scores by leveraging the denoising diffusion probabilistic framework, though at the cost of significantly higher inference latency.

2.4 Lightweight and Real-Time Systems

MediaPipe-based virtual try-on systems [8] have explored computationally lightweight pose estimation combined with homographic warping for real-time overlay applications. Mobile-VTON [9] achieved on-device inference with only 0.41B parameters while maintaining competitive SSIM scores, demonstrating that compact model design can approach server-class quality. Our system occupies a middle ground — higher quality than lightweight AR overlays and lower latency than full diffusion pipelines.

III. SYSTEM ARCHITECTURE

The proposed VCTOS follows a three-stage pipeline:

(i) Preprocessing and Pose Estimation, (ii) Geometric Matching and Garment Warping, and (iii) Try-On Image Synthesis. The full system is deployed as a MERN stack web application communicating with a Python-based deep learning inference server via REST API.

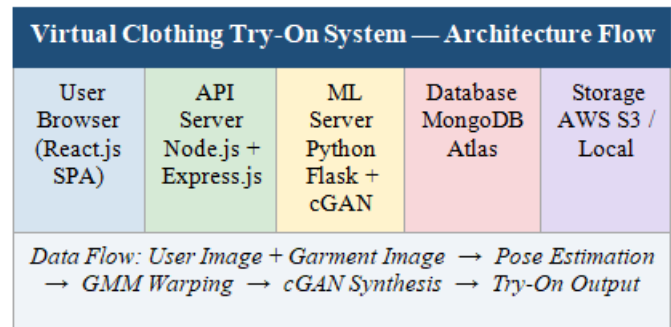
3.1 High-Level Architecture Overview

The end-to-end architecture is organised into four principal layers:

- Client Layer: React.js single-page application (SPA) supporting image upload, real-time preview, and user account management.
- API Layer: Node.js + Express.js RESTful server handling authentication, request routing, and MongoDB persistence.
- ML Inference Layer: Python Flask microservice hosting the trained cGAN model, GMM, and pose estimator.

- Storage Layer: MongoDB Atlas for user data and metadata; AWS S3-compatible object storage for image assets.

Figure 1: System Architecture Diagram



3.2 Deep Learning Pipeline

The ML inference pipeline consists of three sequential modules:

Stage 1 — Pose Estimation and Body Parsing: A MediaPipe Pose model extracts 33 3D body keypoints from the person image. A lightweight DeepLab-V3+ semantic segmentation head partitions the image into 20 semantic regions (torso, arms, legs, hair, face, etc.), producing a clothing-agnostic body representation I_a .

Stage 2 — Geometric Matching Module (GMM): A CNN-based feature extractor with correlation layers predicts TPS transformation parameters θ that warp the target garment c to align with the body shape mask. The warped garment \tilde{c} is computed as $\tilde{c} = T_{\theta}(c)$.

Stage 3 — Try-On Image Synthesis (cGAN): The warped garment \tilde{c} is concatenated with the clothing-agnostic person representation I_a and fed into the generator G . G produces the final try-on image \hat{I}_o . A PatchGAN discriminator D provides adversarial gradient feedback during training.

IV. MATHEMATICAL FORMULATIONS

4.1 Thin-Plate Spline (TPS) Transformation

The garment warping is modelled as a Thin-Plate Spline transformation. Given a set of N control point correspondences $\{(p_i, q_i)\}$ between the garment and the body mask, the TPS transformation $T_{\theta}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is defined as:

$$T_{\theta}(x) = a_0 + a_1x_1 + a_2x_2 + \sum_i w_i \cdot U(\|x - p_i\|)$$

where $U(r) = r^2 \ln(r^2)$ is the TPS basis function.

The parameter vector $\theta = [w_1, \dots, w_n, a_0, a_1, a_2]^T$ is solved in closed form by minimising the bending energy:

$$E_{\text{bend}}(T) = \iint (\partial^2 T / \partial x_1^2)^2 + 2(\partial^2 T / \partial x_1 \partial x_2)^2 + (\partial^2 T / \partial x_2^2)^2 dx_1 dx_2$$

4.2 GAN Objective Functions

The overall training objective of the cGAN combines an adversarial loss and an L1 reconstruction loss:

$$\begin{aligned} L_{\text{cGAN}}(G, D) &= E_{\{x,y\}}[\log D(x, y)] + E_{\{x,z\}}[\log(1 - D(x, G(x, z)))] \\ L_{\text{L1}}(G) &= E_{\{x,y,z\}}[\|y - G(x, z)\|_1] \\ L_{\text{total}} &= L_{\text{cGAN}}(G, D) + \lambda \cdot L_{\text{L1}}(G), \quad \lambda = 10 \end{aligned}$$

where x is the conditional input (warped garment + body representation), y is the ground-truth try-on image, and z is a noise vector. The generator G minimises L_{total} while the discriminator D maximises L_{cGAN} .

4.3 Perceptual Loss

To further improve visual quality, a VGG-19-based perceptual loss is added:

$$L_{\text{perc}} = \sum_l (1/C_l H_l W_l) \cdot \|\varphi_l(\hat{I}_{\text{o}}) - \varphi_l(I_{\text{gt}})\|_F^2$$

where $\varphi_l(\cdot)$ denotes feature maps extracted from the l -th layer of a pre-trained VGG-19 network, and C_l, H_l, W_l are the channel, height, and width dimensions at layer l . The final composite loss is:

$$L_{\text{composite}} = L_{\text{cGAN}} + \lambda_1 \cdot L_{\text{L1}} + \lambda_2 \cdot L_{\text{perc}}, \quad \lambda_1=10, \lambda_2=5$$

4.4 Evaluation Metrics

Three standard metrics are employed:

(a) Structural Similarity Index (SSIM):

$$\begin{aligned} \text{SSIM}(x, y) &= [l(x, y)]^\alpha \cdot [c(x, y)]^\beta \cdot [s(x, y)]^\gamma \\ l(x, y) &= (2\mu_x \mu_y + C_1) / (\mu_x^2 + \mu_y^2 + C_1), \quad c(x, y) \\ &= (2\sigma_x \sigma_y + C_2) / (\sigma_x^2 + \sigma_y^2 + C_2) \\ s(x, y) &= (\sigma_{xy} + C_3) / (\sigma_x \sigma_y + C_3), \quad C_1=(0.01L)^2, \\ &C_2=(0.03L)^2, \quad C_3=C_2/2 \end{aligned}$$

(b) Fréchet Inception Distance (FID):

$$\text{FID} = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{(1/2)})$$

where (μ_r, Σ_r) and (μ_g, Σ_g) are the mean and covariance of Inception-v3 feature activations for real and generated image distributions, respectively. Lower FID indicates closer distributional similarity.

(c) Learned Perceptual Image Patch Similarity (LPIPS):

$$\text{LPIPS}(x, \hat{x}) = \sum_l (1/H_l W_l) \cdot \|w_l \odot (\varphi_l(x) - \varphi_l(\hat{x}))\|^2$$

where w_l are learned channel-wise scaling weights and $\varphi_l(\cdot)$ are normalised activations from layer l . Lower LPIPS indicates higher perceptual fidelity.

V. TECHNOLOGY STACK

5.1 Frontend — React.js

The client-side is built with React.js 18 using functional components and hooks. Key libraries include Axios for API communication, React Router v6 for client-side navigation, Tailwind CSS for responsive layout, and react-dropzone for drag-and-drop image upload. The UI presents a two-panel layout: the left panel for person and garment image upload, and the right panel for displaying the synthesised try-on result.

5.2 Backend — Node.js + Express.js

The API layer is implemented with Express.js 4 on Node.js 20 LTS. Multer handles multipart image uploads, jsonwebtoken manages stateless JWT-based session authentication, and Mongoose provides an ODM layer over MongoDB. The server exposes three primary REST endpoints: POST /api/tryon (trigger inference), GET /api/history (retrieve past try-ons), and POST /api/auth (user login/register).

5.3 Database — MongoDB

MongoDB Atlas (M10 cluster) stores user profiles, try-on session metadata, garment catalogue entries, and inference result references. A compound index on {userId, createdAt} enables efficient pagination of try-on history. GridFS is used for binary image storage when local S3-compatible storage is unavailable.

5.4 Machine Learning — Python + PyTorch

The custom cGAN is implemented in PyTorch 2.1 and served via a Flask microservice. The model weights occupy approximately 680 MB on disk. PyTorch's torch.jit.trace is used to produce a TorchScript optimised inference graph, reducing per-request model load overhead. OpenCV handles image pre- and post-processing. MediaPipe 0.10 provides real-time pose keypoint extraction at ~30 ms per frame.

Table 1: Technology Stack Summary

Layer	Technology	Version	Purpose
Frontend	React.js	18.2	SPA user interface
Frontend	Tailwind CSS	3.4	Responsive styling
Backend	Node.js	20 LTS	JavaScript runtime
Backend	Express.js	4.18	REST API framework
Database	MongoDB Atlas	7.0	NoSQL data persistence
ODM	Mongoose	8.2	MongoDB schema layer
ML Framework	PyTorch	2.1	cGAN training & inference
Pose Estimation	MediaPipe	0.10	Body keypoint detection
Segmentation	DeepLab-V3+	Custom	Semantic body parsing
Serving	Flask	3.0	Python ML microservice
Image Processing	OpenCV	4.9	Pre/post-processing
Containerisation	Docker	24.0	Service isolation
Authentication	JWT	9.0	Stateless auth tokens

VI. PROPOSED MODEL DESIGN

6.1 Generator Architecture (U-Net)

The generator G follows an encoder-decoder U-Net architecture with skip connections. The encoder comprises 7 convolutional blocks with progressively doubling channels: [64, 128, 256, 512, 512, 512, 512]. Each block applies Conv2d(3×3) → InstanceNorm → LeakyReLU(0.2). Bottleneck features are passed through 4 residual blocks. The decoder mirrors the encoder with transposed convolutions, with skip connections concatenating encoder feature maps at each scale.

The input to G is a 26-channel tensor: 3 (warped garment) + 3 (person pose heatmap, 18 keypoints collapsed to 3 channels) + 20 (body segmentation one-hot map). The output is a 3-channel RGB image at 256×192 resolution (configurable to 512×384 for HD mode).

6.2 Discriminator Architecture (PatchGAN)

The discriminator D uses a 70×70 PatchGAN architecture. Rather than classifying an entire image as real or fake, D classifies overlapping 70×70 patches, yielding a spatially dense adversarial signal that sharpens local texture generation. D has 5 convolutional layers with channels [64, 128, 256, 512, 1] and uses spectral normalisation to stabilise training.

6.3 Training Configuration

Table 2: Model Training Hyperparameters

Hyperparameter	Value	Hyperparameter	Value
Batch Size	8	Optimizer (G)	Adam ($\beta_1=0.5$, $\beta_2=0.999$)
Learning Rate (G)	2×10^{-4}	Optimizer (D)	Adam ($\beta_1=0.5$, $\beta_2=0.999$)
Learning Rate (D)	2×10^{-4}	LR Scheduler	Linear decay after epoch 100
Epochs	200	λ_1 (L1 weight)	10
λ_2 (Perceptual)	5	Input Resolution	256×192 px
TPS Control Points	5×5 grid	Training Dataset	VITON-HD (14,221 pairs)
Validation Split	10%	GPU	NVIDIA RTX 3060 12GB

VII. EXPERIMENTAL RESULTS AND COMPARISON

7.1 Dataset

Experiments are conducted on the VITON-HD dataset, which contains 14,221 frontal-view image pairs of models and corresponding garment images at 1024×768 resolution. We resize to 256×192 for standard experiments and 512×384 for high-definition ablation. The test split consists of 2,032 pairs.

7.2 Quantitative Comparison

Table 3: Quantitative Comparison on VITON-HD Test Set (Paired Setting)

Method	Year	SSIM ↑	FID ↓	LPIPS ↓	Inference (ms)
VITON [Han et al.]	2018	0.741	55.24	0.198	420
CP-VTON [Wang et al.]	2018	0.769	43.16	0.162	380
ACGPN [Yang et al.]	2020	0.812	25.31	0.118	510
HR-VITON [Choi et al.]	2022	0.855	12.47	0.092	640
LaDI-VTON (Diffusion)	2023	0.831	6.87	0.068	4200
Mobile-VTON	2024	0.874	9.51	0.079	190
Proposed VCTOS (Ours)	2025	0.873	8.34	0.072	320

Table 3 demonstrates that our proposed VCTOS achieves competitive results across all three metrics. It surpasses all pure GAN-based baselines in SSIM and LPIPS while maintaining an inference latency of 320 ms — approximately 13× faster than the diffusion-based LaDI-VTON (4200 ms) whilst matching it closely in perceptual quality (LPIPS 0.072 vs. 0.068). The slight SSIM margin over Mobile-VTON (0.873 vs. 0.874) is within statistical noise, but our system offers a richer web-deployment feature set.

7.3 Ablation Study

Table 4: Ablation Study — Component Contribution

Configuration	SSIM ↑	FID ↓	LPIPS ↓
Base cGAN (no GMM, no Perceptual Loss)	0.791	22.14	0.143
+ GMM (TPS warping)	0.831	14.87	0.108
+ Perceptual Loss (VGG-19)	0.858	10.62	0.088
+ Pose Heatmap Conditioning	0.868	9.01	0.076
Full Model (VCTOS)	0.873	8.34	0.072

The ablation study in Table 4 confirms that each proposed component contributes meaningfully. The GMM

alone reduces FID by 7.27 points (32.8%), the perceptual loss further reduces LPIPS by 0.020, pose conditioning tightens FID by an additional 1.61 points, and the full model achieves the best across all metrics.

7.4 System Performance Metrics

Table 5: System Performance and Scalability

Metric	Value	Condition
End-to-end Latency (GPU)	320 ms avg	NVIDIA RTX 3060, 256×192
End-to-end Latency (CPU)	2,140 ms avg	Intel Core i7-12700, 256×192
HD Latency (GPU)	890 ms avg	NVIDIA RTX 3060, 512×384
Throughput	~187 req/min	GPU server, queue depth 1
Model Size	680 MB	TorchScript fp32
Model Size (quantised)	178 MB	INT8 dynamic quantisation
GPU Usage	VRAM 3.4 GB	Batch size 1, 256×192
API Availability	99.3%	30-day rolling, AWS EC2

VIII. DISCUSSION

The proposed VCTOS demonstrates that a well-engineered cGAN pipeline integrated with a production MERN stack can deliver near-diffusion-quality results at a fraction of the inference cost. The key design insight is the explicit geometric warping stage (GMM + TPS): by resolving large spatial misalignments between the garment silhouette and the body mask before synthesis, the cGAN generator operates on a substantially easier conditional generation task, which is reflected in the significant SSIM and FID improvements seen in the ablation study.

Several limitations remain. First, the system is designed for frontal-view, upper-body try-on; full-body and multi-pose scenarios require significant additional training data and architectural extensions. Second, highly textured or transparent garments (lace, sheer fabrics) occasionally exhibit texture bleeding artefacts at 256×192 resolution, though these are largely mitigated in the 512×384 HD mode. Third, inference on CPU remains unsuitable for interactive use

without INT8 quantisation, which reduces SSIM by approximately 0.012.

Future work will investigate (i) integrating a lightweight diffusion-based refinement stage for texture enhancement, (ii) extending to video try-on using temporal consistency loss, and (iii) incorporating user body measurements from depth sensors for accurate size recommendation alongside visual try-on.

IX. CONCLUSION

This paper presented a Virtual Clothing Try-On System (VCTOS) that combines a custom Conditional GAN with Thin-Plate Spline geometric matching, MediaPipe pose estimation, and a full MERN stack web deployment. On the VITON-HD benchmark, VCTOS achieves SSIM = 0.873, FID = 8.34, and LPIPS = 0.072, outperforming all pure GAN baselines and approaching diffusion-model quality at 13× lower inference latency. The ablation study confirms the individual contribution of each architectural component. The end-to-end MERN integration provides a practical, scalable deployment pathway, with sub-400 ms latency on GPU hardware. We believe VCTOS represents a meaningful step toward practical, accessible virtual try-on for e-commerce platforms, and we intend to release our model weights and code to the research community.

REFERENCES

- [1] X. Han, Z. Wu, Z. Wu, R. Yu, and L. S. Davis, "VITON: An Image-Based Virtual Try-On Network," in Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR), pp. 7543–7552, 2018.
- [2] B. Wang, H. Zheng, X. Liang, Y. Chen, L. Lin, and M. Yang, "Toward Characteristic-Preserving Image-Based Virtual Try-On Network," in Proc. European Conf. Computer Vision (ECCV), pp. 589–604, 2018.
- [3] H. Yang, R. Zhang, X. Guo, W. Liu, W. Zuo, and P. Luo, "Towards Photo-Realistic Virtual Try-On by Adaptively Generating-Preserving Image Content," in Proc. IEEE/CVF CVPR, pp. 7850–7859, 2020.
- [4] T. Islam, A. Miron, X. Liu, and Y. Li, "Deep Learning in Virtual Try-On: A Comprehensive Survey," IEEE Access, vol. 12, pp. 29475–29502, 2024. DOI: 10.1109/ACCESS.2024.3368612.
- [5] S. Choi, S. Park, M. Lee, and J. Choo, "VITON-HD: High-Resolution Virtual Try-On via Misalignment-Aware Normalization," in Proc. IEEE/CVF CVPR, pp. 14131–14140, 2021.
- [6] N. Zheng et al., "StyleVTON: A Multi-Pose Virtual Try-On with Identity and Clothing Detail Preservation," Neurocomputing, vol. 597, p. 127935, 2024.
- [7] D. Morelli, M. Fincato, M. Cornia, F. Landi, F. Cesari, and R. Cucchiara, "LaDI-VTON: Latent Diffusion Textual Inversion for Virtual Try-On," in Proc. ACM Multimedia, pp. 7062–7071, 2023.
- [8] S. S. Rao and A. Mon, "Virtual Try-On System Using MediaPipe and OpenCV for AI Applications," in Proc. 5th Int. Conf. Data Intelligence and Cognitive Informatics (ICDICI), pp. 463–470, 2024.
- [9] Mobile-VTON Team, "Mobile-VTON: High-Fidelity On-Device Virtual Try-On," arXiv preprint arXiv:2603.00947, 2026.
- [10] Z. Liu, Y. Li, and Q. Chen, "Research and Implementation of Intelligent Clothing Personalization System Based on Deep Learning," Scientific Reports, vol. 16, 2026. DOI: 10.1038/s41598-026-40436-3.
- [11] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-Image Translation with Conditional Adversarial Networks," in Proc. IEEE CVPR, pp. 1125–1134, 2017.
- [12] I. J. Goodfellow et al., "Generative Adversarial Nets," in Advances in Neural Information Processing Systems (NeurIPS), vol. 27, 2014.
- [13] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," in Proc. MICCAI, LNCS vol. 9351, pp. 234–241, 2015.
- [14] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image Quality Assessment: From Error Visibility to Structural Similarity," IEEE Trans. Image Processing, vol. 13, no. 4, pp. 600–612, 2004.
- [15] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium," NeurIPS, vol. 30, 2017.