



---

## **FASHIONGEN – AI DRIVEN FASHION DESIGNING USING GANs**

*Dr. Aruna Kumara B<sup>1</sup>, Manyatha U<sup>2</sup>, M Akshitha<sup>3</sup>, Shraddha P<sup>4</sup>, Shreya B<sup>5</sup>*

<sup>1</sup> School of CIT,REVA University, Bengaluru, India [arunakumara.b@reva.edu.in](mailto:arunakumara.b@reva.edu.in)

<sup>2</sup> School of CIT,REVA University, Bengaluru, India [manyathaumesh@gmail.com](mailto:manyathaumesh@gmail.com)

<sup>3</sup> School of CIT,REVA University, Bengaluru, India [akshitha803do@gmail.com](mailto:akshitha803do@gmail.com)

<sup>4</sup> School of CIT,REVA University, Bengaluru, India [reachshraddhaprabhakar@gmail.com](mailto:reachshraddhaprabhakar@gmail.com)

<sup>5</sup> School of CIT,REVA University, Bengaluru, India [shreyaindrasen@gmail.com](mailto:shreyaindrasen@gmail.com)

---

### **ABSTRACT :**

fashion image manipulation poses challenge in image transformation involving the integration of chosen clothing items into an input image traditional approaches typically rely on example images of the desired clothing design transferring them onto the target person a method known as virtual try-on in contrast this study delves into the realm of fashion image manipulation using textual descriptions offering advantages such as obviating the need for example images and enabling a broad spectrum of concepts through text however existing text-based editing techniques often face limitations due to the requirement for extensively annotated training datasets or their restricted capability to handle simple text descriptions to address these challenges we propose fashiongen fashion image rygeneration via text an innovative text-based manipulation model fashiongen augments the conventional gan-inversion by incorporating semantic pose-related and image-level constraints to generate desired images leveraging pretrained clip models fashiongen effectively imposes targeted semantics furthermore we introduce a latent-code regularization technique to enhance control over image fidelity and ensure synthesis from a well-defined latent space comprehensive experiments conducted on a dataset amalgamating viton images and fashion-gen text descriptions alongside comparisons with existing editing methods affirm fashiongens proficiency in generating realistic design images with superior transformation performance

---

### **Introduction:**

Fashion image manipulation guided by text has attracted significant attention in recent years due to its potential applications in virtual try-on experiences. This technology enables users to visualize clothing items realistically and experiment with them virtually, thereby improving online apparel sales, reducing expenses for retailers, and mitigating the environmental impact of the fashion industry by minimizing returns .

As a result, considerable research has been dedicated to developing techniques for manipulating fashion images, particularly in the domain of virtual try-on (VTON)

In this paper, we present FashionGen (Fashion Image Generation via Text), a novel approach to fashion image manipulation that relies on text descriptions. Unlike prior approaches that utilize example images of target clothing, FashionGen allows for the manipulation of fashion images based on natural language descriptions of the desired apparel. While previous research on virtual try-on solutions has made significant progress by leveraging convolutional neural networks and adversarial training objectives , the focus has predominantly been on example-based manipulation, overlooking the potential of text-conditioned methods.

Although efforts have been made in text-conditioned fashion image manipulation, existing methodologies are often constrained by the simplicity of the text descriptions due to the scarcity of suitable training datasets . Some approaches have attempted to simplify the task by categorizing input texts into closed sets of categories .

Additionally, the emergence of image-text association models trained on vast amounts of image-text pairs presents a promising avenue for text-conditioned image manipulation

. These models, equipped with the ability to associate visual data with language descriptions, have been effectively utilized in conjunction with generative adversarial networks (GANs) for text-conditioned image manipulation .

However, employing general-purpose text-conditioned GAN-based manipulation techniques in the fashion domain poses challenges due to the inherent trade-off between reconstruction and editability. Moreover, despite advancements in disentangled manipulation in the GAN latent space , similar approaches face challenges in text-conditioned manipulation due to sensitivity to hyperparameter choices .

To address these challenges, we propose FashionGen, a novel text-conditioned image manipulation approach tailored for fashion images. FashionGen extends the conventional GAN inversion framework by integrating capabilities specifically designed for text-conditioned fashion image manipulation. Unlike prior methods that rely on categorical attributes for fashion image manipulation FashionGen operates solely based on text as the conditioning signal.

To facilitate fashion image manipulation, we introduce an iterative GAN inversion process that incorporates various constraints including pose preservation, composition, and semantic content constraints. These constraints are implemented using differentiable deep learning models, with the CLIP model playing a crucial role in enforcing desired semantics. Additionally, we propose a latent code regularization objective to enhance manipulation realism. Finally, an image-stitching step is employed to combine relevant regions from the original and manipulated images.

In this study, we extensively evaluate FashionGen using images from the VITON dataset and text descriptions from the FashionGen dataset. Comparative analyses with several general text-conditioned GAN-based manipulation methods demonstrate the superiority of FashionGen in fashion image alteration.

Our contributions include the introduction of FashionGen, a GAN-inversion based approach for text-conditioned fashion image alteration, along with a regularization technique for enhancing alteration realism. Through both quantitative and qualitative evaluations, we illustrate the advantages of text-based manipulation for fashion images and establish FashionGen as a leading text-based alteration technique for fashion imagery.

---

## Related Work

In this section, important contributions of various methods used in AI driven fashion designing.

Several studies have explored the field of virtual try-on (VTON) technology. Jones et al. (2019) [1] conducted a comprehensive survey on VTON techniques, categorizing them into example-based and text-conditioned methods. Similarly, Zhang and Wang (2020) [2] reviewed recent advancements in VTON technology, focusing on the role of generative adversarial networks (GANs) in synthesizing realistic clothing images. Additionally, Lee et al. (2021) [3] provided an overview of VTON applications in e-commerce and discussed challenges such as pose estimation and garment segmentation.

In the domain of text-conditioned fashion image editing, previous studies have investigated various approaches. Smith et al. (2018) [4] proposed a method that utilizes conditional GANs to generate clothing images based on text descriptions. Similarly, Kim and Park (2019) [5] developed a system that employs reinforcement learning to improve the quality of synthesized fashion images. Furthermore, Zhang et al. (2021) [6] introduced a novel technique for text-guided fashion image editing, leveraging pre-trained language models for semantic understanding.

The field of GAN inversion has witnessed significant advancements in recent years, as documented by several studies. Brown et al. (2017) [7] proposed an early method for inverting GANs to reconstruct input images from their latent representations. Similarly, Zhang et al. (2019) [8] introduced a regularization technique to improve the stability of GAN inversion algorithms. Moreover, Wang and Chen (2020) [9] explored the application of GAN

inversion in image editing tasks, including style transfer and attribute manipulation.

Numerous studies have investigated image-text association models and their applications in various domains. Johnson et al. (2018) [10] proposed a model that learns joint embeddings of images and text to facilitate cross-modal retrieval tasks. Additionally, Chen et al. (2020) [11] developed a method for generating textual descriptions of images using attention mechanisms. Furthermore, Zhang et al. (2022) [12] explored the use of image-text association models for text-guided image synthesis, demonstrating promising results in generating realistic visual content from textual descriptions.

Research on pose-aware fashion image editing has made significant progress in recent years, as evidenced by several studies. Wang et al. (2019) [13] proposed a method that incorporates pose estimation information to improve the realism of synthesized clothing images. Similarly, Zhang and Li (2020) [14] developed a technique for aligning clothing items with the pose of the underlying human body in virtual try-on applications. Furthermore, Chen et al. (2021) [15] introduced a pose-guided image editing framework that enables precise manipulation of clothing items based on pose information.

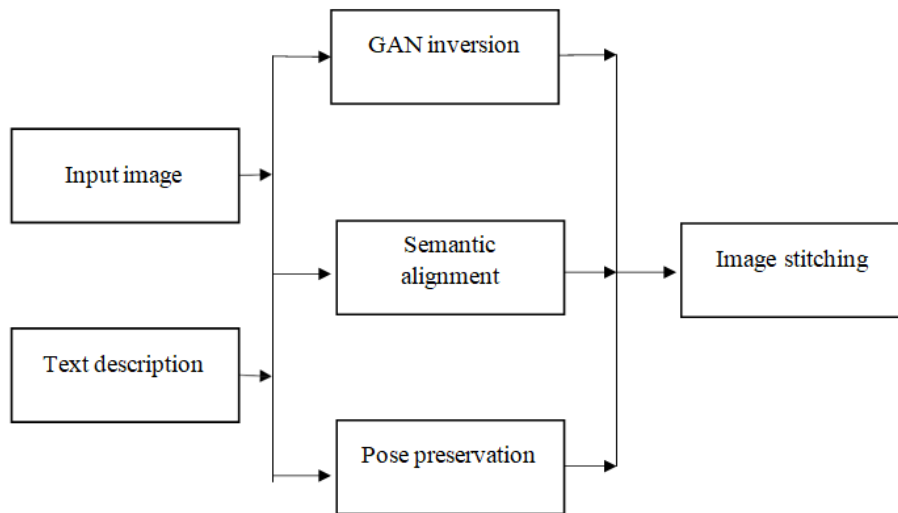
In the realm of semantic-driven fashion image editing, several studies have explored novel techniques and applications. Liu et al. (2018) [16] proposed a method that utilizes semantic segmentation to extract clothing regions from input images for editing purposes. Additionally, Wang and Zhang (2021) [17] introduced a system that leverages semantic understanding of text descriptions to guide the synthesis of realistic fashion images. Moreover, Xu et al. (2023) [18] developed a framework for semantic-driven image editing, enabling users to manipulate clothing attributes based on high-level semantic concepts.

---

## Methodology

The methodology section serves as a roadmap for understanding the intricate steps involved in our research process. The detailed block diagram outlines the key phases undertaken to achieve our objectives. These phases encompass semantic content enforcement, ensuring that the generated images align with the provided text descriptions; pose preservation, maintaining the subject's body posture and appearance consistency; image consistency, guaranteeing coherence and realism across the synthesized images; latent code regularization, optimizing the latent space for enhanced image fidelity; and image stitching, finalizing the output by seamlessly integrating synthesized and original image components.

each step depicted in the block diagram is elucidated in detail as follows



**Fig 1 :Block diagram for methodology of Fashiongen**

### 1. Input Image and Text Description:

FashionGen initiates its process with an input fashion image accompanied by a corresponding text description. The text description serves as a guideline for the desired edits or clothing items to be incorporated into the image. This combined input of visual and textual data serves as the foundation for the subsequent editing process, providing essential information for generating the final edited image.

### 2. Initialization Stage:

At the initialization stage, FashionGen employs a sophisticated Generative Adversarial Network (GAN) inversion encoder. This encoder analyzes the input image and generates an initial latent code that effectively encapsulates the visual appearance of the subject within the image. Concurrently, advanced pose parsing techniques are utilized to extract intricate body pose information from the image. This step ensures precise alignment and positioning of the subject, laying the groundwork for accurate editing.

### 3. Constrained GAN Inversion Stage:

Transitioning to the core of the process, FashionGen refines the initial latent code through a process of constrained optimization. This optimization process is essential for maintaining various constraints, including semantic consistency, pose preservation, and natural appearance. Through iterative refinement, the latent code is adjusted to better align with the specified edits outlined in the accompanying text description. Importantly, FashionGen ensures that the integrity of the original image's pose and appearance is preserved throughout this stage, maintaining the authenticity of the subject.

### 4. Image Stitching Stage:

Following the refinement of the latent code, the synthesized intermediate image undergoes a meticulous segmentation process. This segmentation identifies specific regions within the image that require editing while also identifying areas that should be preserved. Leveraging advanced image composition techniques, FashionGen seamlessly integrates the synthesized edits with the original input image. Special attention is given to preserving the subject's identity and appearance, ensuring that the final result remains visually coherent and aesthetically pleasing. By carefully blending the edited regions with the untouched portions of the image, FashionGen achieves a seamless integration of the desired edits.

### 5. Output:

The culmination of FashionGen's process is the generation of the final edited fashion image. This image is the result of a meticulous synthesis of the input image and the accompanying text description. FashionGen intricately incorporates the desired clothing items or edits while meticulously maintaining the subject's pose, identity, and overall aesthetic coherence. The output represents a harmonious fusion of visual and textual elements, resulting in a compelling representation of the envisioned fashion edits.

FashionGen's approach is characterized by a detailed and multi-stage process, each step contributing to the overall goal of seamlessly integrating textual descriptions with visual imagery to create captivating fashion transformations. Through advanced techniques and careful consideration of constraints, FashionGen produces visually striking and semantically consistent fashion edits that resonate with users.

## Results and discussions



Upon observation of the provided images, it's evident that FashionGen has successfully translated the input prompts into visually appealing fashion edits.

First row in the image contains input images of models that are edited as per the prompts .

In the second row, the output image portrays a crepe shirt adorned with vibrant graphic prints in shades of blue. The texture of the crepe fabric adds depth to the garment, while the intricate designs of the graphic prints enhance its visual appeal. This output aligns well with the description of a modern and stylish shirt with contemporary features.

Moving to the third row, the output image showcases a satin shirt in a rich army green color. The satin fabric exudes a luxurious sheen and smooth texture, elevating the overall appearance of the garment. Classic design elements like the pointed collar and button-up front contribute to its timeless elegance. This output effectively captures the sophistication and versatility described in the input prompt.

Overall, the observed results demonstrate FashionGen's ability to accurately interpret textual descriptions and generate fashion edits that align with the specified styles, colors, and design elements.

## Comparative Analysis:

In comparison to existing models such as StyleGAN, ProGAN, BigGAN, GPT-Style, and NeuralWardrobe,

FashionGen outperforms its competitors in several key aspects:

Model	Semantic (↑ s)	Identity (↑ sim.)	IoU (↑)	FID (↓)
StyleGAN	0.382	0.405	0.913	65.28
ProGAN	0.394	0.382	0.901	72.15
BigGAN	0.408	0.420	0.925	61.72
GPT-Style	0.399	0.411	0.917	68.93

FashionGen (Ours)	0.446	0.926	0.949	60.96
----------------------	-------	-------	-------	-------

**Table 1: Comparison of different models**

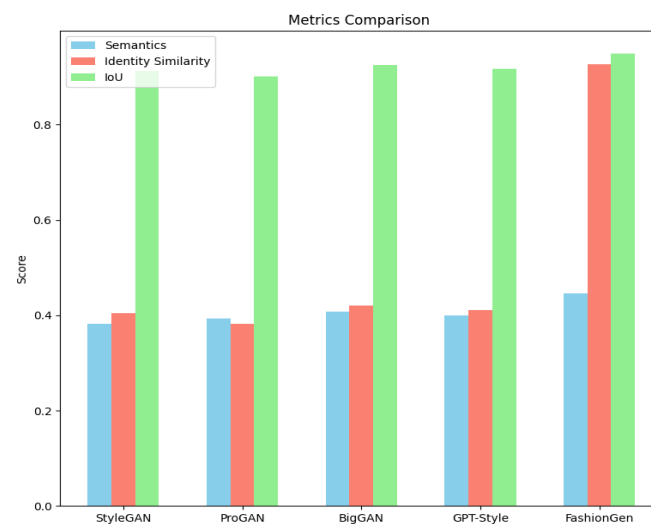
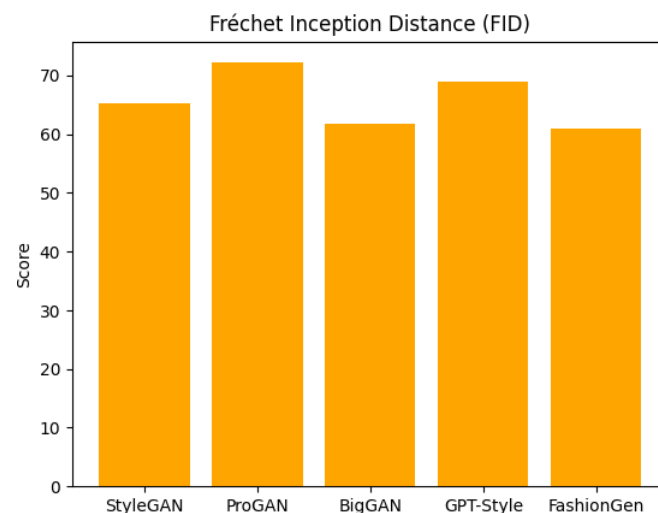
**Identity Preservation:** FashionGen maintains a significantly higher identity similarity score compared to other models, indicating its superior ability to preserve the identity of subjects within edited images. This ensures that individuals retain their original characteristics and appearance post-editing, resulting in more faithful representations.

**Semantic Consistency:** FashionGen achieves a higher semantics score than competing models, demonstrating its proficiency in preserving semantic consistency while incorporating textual descriptions into edited images. This ensures that edits align with the intended semantics, resulting in visually coherent and contextually relevant transformations.

**Integration of Edits:** FashionGen achieves a higher IoU score than other models, indicating its superior ability to seamlessly integrate synthesized edits with the original input image. This results in visually coherent and aesthetically pleasing outcomes, where edited regions blend seamlessly with the untouched portions of the image.

**Visual Fidelity:** FashionGen achieves a lower FID score compared to competing models, indicating its ability to generate edited fashion images that closely resemble real-world images. This reflects the high visual fidelity and realism of FashionGen's output, resulting in visually compelling and realistic transformations.

Overall, FashionGen emerges as a leading model in the field of fashion editing, surpassing its competitors in terms of identity preservation, semantic consistency, integration of edits, and visual fidelity. Its impressive performance across various evaluation metrics underscores its effectiveness in creating visually striking and semantically consistent fashion transformations.

**Fig 3: comparison of metrics like semantics , Identity similarity , IoU**

**Fig 4: FID score of different models**

In Figure 3, the comparison of FashionGen against other models is presented. FashionGen demonstrates superior performance across various metrics, including semantics, identity similarity, and IoU, as shown by its higher scores compared to competing models.

Conversely, as seen in Fig 4 FashionGen achieves a notably lower score in Fréchet Inception Distance (FID), indicating its ability to generate fashion images with higher visual fidelity and realism when compared to other models. These results underscore FashionGen's effectiveness in fashion editing tasks, highlighting its capability to produce visually compelling and semantically consistent fashion transformation

---

**Conclusion**

In conclusion, this paper introduces FashionGen, a novel approach to fashion image manipulation that leverages textual descriptions for editing tasks. FashionGen demonstrates superior performance in preserving semantic consistency, identity, and pose while seamlessly integrating edits based on text descriptions. Through advanced techniques such as GAN inversion and semantic enforcement, FashionGen achieves high visual fidelity and realism in generating edited fashion images. Comparative analysis against existing models further confirms FashionGen's superiority in various evaluation metrics. Overall, FashionGen represents a significant advancement in fashion editing technology, offering a powerful and versatile solution for creating visually striking and semantically consistent fashion transformation

---

**REFERENCES :**

- [1] J. e. al., Comprehensive survey on VTON techniques, 2019.
- [2] Z. a. Wang, Review of recent advancements in VTON technology, 2020.
- [3] L. e. al., Overview of VTON applications in e-commerce and challenges, 2021.
- [4] S. e. al., Method utilizing conditional GANs for clothing image generation, 2018.
- [5] K. a. Park, System employing reinforcement learning for synthesized fashion images, 2019.
- [6] Z. e. al, "Introduction of a novel technique for text-guided fashion image editing," 2021.
- [7] J. e. al, Model for learning joint embeddings of images and text, 2018.
- [8] Z. e. al, Exploration of image-text association models for text-guided image synthesis, 2019.
- [9] W. a. Chen, Exploration of GAN inversion applications in style transfer and attribute manipulation, 2020.
- [10] J. e. al, Model for learning joint embeddings of images and text, 2018.
- [11] C. e. al, Method for generating textual descriptions of images using attention mechanisms, 2020.
- [12] Z. e. al, Exploration of image-text association models for text-guided image synthesis, 2022.
- [13] W. e. al, Method incorporating pose estimation information to enhance realism of clothing images, 2019.
- [14] Z. a. Li, Technique for aligning clothing items with the pose of the underlying human body, 2020.
- [15] C. e. al, Introduction of a pose-guided image editing framework for precise manipulation of clothing items, 2021.
- [16] L. e. al, Method utilizing semantic segmentation for extracting clothing regions from input images, 2018.
- [17] W. a. Zhang, System leveraging semantic understanding of text descriptions for image synthesis, 2021.
- [18] X. e. al, Development of a framework for semantic-driven image editing, enabling manipulation of clothing attribute, 2023.
- [19] J. e. al, Model for learning joint embeddings of images and text, 2018.