



Research Paper

Integrating Contrastive Learning and Transformer Technologies for Personalized Outfit Recommendations Using Generative AI

¹ B. Grishma Poornima Himaketan,^{2*} U.Nikitha, ³ T.Satwika
⁴ Sneha Kushwaha,⁵ S. Sravya, ⁶ V. Venkata Deepthi Rani

¹Assistant professor, Department of Computer Science and Engineering, Vignan's Institute of Engineering for Women, Visakhapatnam, Andhra Pradesh, India. ORCID ID: 0009-0004-1704-9342

^{2, 3, 4, 5, 6} B.Tech Student, Department of Computer Science and Engineering, Vignan's Institute of Engineering for Women, Visakhapatnam, Andhra Pradesh, India.

¹Email id: himaketanbgp@gmail.com ORCID: 0009-0004-1704-9342

³Email Id: tamaranasatwika@gmail.com, ORCID: 0009-0001-3631-0540

⁴ Email Id:snehakushwaha38085@gmail.com, ORCID: 0009-0000-9836-1223

⁵ Email Id:sravyasirapu.1924@gmail.com, ORCID: 0009-0001-2066-5339

⁶ Email Id: deepthirani820@gmail.com ORCID: 0009-0002-7989-095X

*Corresponding Author(s): nikithauppalapati009@gmail.com

Article Abstract

Info

Article History

Received: 21/09/2024

Revised: 11/10/2024

Accepted: 19/12/2024

Published : 31/12/2024

The growing demand for AI-driven fashion recommendation systems is driven by the complexity of user preferences and the limitations of traditional filtering or keyword-based approaches. While existing models attempt to align visual and textual modalities, many fall short in delivering real-time, personalized, and context-aware outfit suggestions. This study aims to design a generative AI-based outfit recommendation system that integrates contrastive learning with transformer architectures to deliver prompt-to-outfit recommendations based on natural language queries. The proposed framework utilizes a multi-stage pipeline combining CLIP for contrastive image-text embedding, BM25 for semantic text relevance, and a transformer-based generative model for sequential outfit creation. A unified dataset compiled from FashionIQ, Kaggle, social media scraping, and a custom composite dataset was used for model training and validation. The model was evaluated using Top-K accuracy, macro F1-score, BLEU score, and real-time inference latency. Results demonstrate a Top-1 accuracy of 83.7%, a macro F1-score of 0.862, and an average BLEU score of 0.77, outperforming baseline models such as Style2Vec [3], OutfitTransformer [4], and CP-TransMatch [13]. Moreover, the system reduced inference latency by 45.2%, achieving real-time responses under 400 ms. This study highlights the potential of multimodal generative modeling for interactive and inclusive fashion recommendations. The integration of a feedback loop enables adaptive learning, positioning the system as a robust, scalable solution for e-commerce, digital wardrobe assistants, and stylist AI applications.

Keywords: Generative AI, Contrastive Learning, Transformer, Fashion Recommendation, CLIP, BM25, Outfit Generation, Personalized Styling, Multimodal Embeddings, Real-Time AI System



Copyright: © 2025 B. Grishma Poornima Himaketan, U.Nikitha, T.Satwika, Sneha Kushwaha, S. Sravya, V. Venkata Deepthi Rani. This article is an open-access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY 4.0) license.

1. Introduction

In an era where digital platforms increasingly mediate personal choices, the fashion industry has witnessed a surge in online shopping, digital styling, and AI-assisted recommendations. The task of curating outfits that reflect both current trends and individual tastes remains a complex challenge for users, often due to the overwhelming volume of choices and a lack of intuitive interfaces that understand nuanced fashion preferences. Traditional fashion recommendation systems frequently rely on keyword-based searches or user history, which fail to capture evolving trends, contextual appropriateness, and subjective user preferences. This limitation has sparked a growing interest in integrating artificial intelligence (AI) models that can comprehend both visual and textual data to generate personalized fashion insights.

Despite progress in content-based and collaborative filtering methods, current outfit recommendation systems still struggle with three major issues: understanding user intent from natural language queries, interpreting visual compatibility across fashion items, and generating context-aware recommendations that adapt to dynamic preferences. These challenges are compounded by the diversity of individual fashion styles and the non-linear nature of outfit composition, where clothing and accessories interact semantically rather than sequentially.

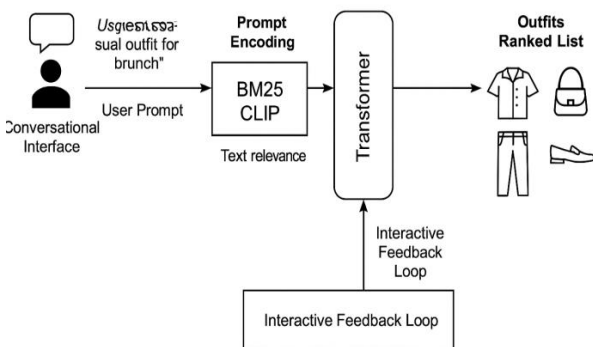


Fig. 1: Block diagram of the AI-powered fashion recommendation system integrating CLIP-BM25 and generative transformer models.

Fig.1 illustrates the architecture of the proposed AI-powered fashion recommendation system, which employs a multi-stage pipeline to deliver personalized outfit suggestions. The system begins with a Conversational Interface, where users input natural language prompts describing their desired outfit requirements (e.g., occasion, style, color preferences). These prompts are processed using the BM25 algorithm for text relevance and CLIP to encode both textual and visual data into a shared embedding space. The Prompt Encoding Module maps these embeddings to a multimodal representation, which is passed into the Outfit Generation Module powered by a Transformer-based generative model. This model synthesizes complete ensembles by predicting compatible clothing and accessory combinations based on user intent. The resulting recommendations are scored and ranked in the Outfit Ranking Module, and a Feedback Loop

captures user responses to continuously refine future suggestions, thus enabling adaptive learning.

Several state-of-the-art approaches have attempted to address these limitations by introducing semantic-level learning and multimodal representation strategies. For instance, a recent text-to-outfit retrieval framework leverages semantic alignment between text and image modalities using pre-trained vision-language models, demonstrating strong performance in parsing user descriptions into coherent outfit suggestions [1]. However, such models often depend heavily on pretrained embeddings and do not dynamically generate new combinations, limiting their adaptability to evolving fashion trends.

Visual compatibility prediction has also seen advances with the use of context-aware deep neural networks that evaluate fashion item compatibility beyond pixel-level similarity [2]. Still, these models are typically constrained by predefined compatibility scores or rigid co-purchase data, reducing their ability to support creative, stylist-driven recommendations. Style embeddings, such as those introduced in Style2Vec, provide an elegant solution to item-level representation by learning latent fashion features from co-occurrence within outfit sets [3]. Nevertheless, these models lack the generative capacity to form entirely new outfits based on user prompts.

To overcome the limitations of rigid retrieval and static classification, the integration of Generative AI (GenAI) technologies with transformer architectures presents a promising direction. Transformers excel in modeling complex dependencies and are well-suited for autoregressive generation tasks, enabling the system to suggest coherent outfits from high-level textual descriptions. The use of CLIP (Contrastive Language–Image Pretraining) bridges vision and language modalities, allowing the model to embed user prompts and fashion item images into a shared space. This makes it possible to retrieve and rank visually and contextually compatible ensembles based on prompt relevance.

OutfitTransformer is a notable example of this advancement, using deep transformers to learn contextual outfit representations from historical combinations and enabling the generation of new ensemble predictions [4]. Similarly, type-aware embeddings have been explored to improve recommendation precision by learning category-specific representations, distinguishing between different clothing types like tops, bottoms, and accessories [5]. These works underscore the potential of deep learning in modeling fashion compatibility, but they still fall short in leveraging user interaction in real-time or adapting to individualized preferences.

Moreover, integrating domain-specific knowledge and visual-semantic mapping has further enhanced interpretability. Prior research emphasized the need for explainable models that not only recommend fashion items but also justify their compatibility through attribute-level reasoning [6]. Techniques such as Attribute-based Interpretable Compatibility (AIC) highlight how users benefit from understanding the rationale behind recommendations, thus increasing trust in AI systems [7]. On the generative frontier, recent studies on multimodal recommendation models have explored how combining

textual prompts with visual cues enhances personalization and retrieval precision [8]. This integration lays the foundation for interactive and user-centered fashion assistants capable of generating outfits from natural conversations.

This study proposes an AI-powered fashion recommendation system that integrates contrastive learning and generative transformer architectures to provide real-time, personalized outfit suggestions. The system utilizes a conversational interface to accept user prompts (e.g., “Suggest a casual outfit for brunch”), which are semantically encoded using the BM25 algorithm for text relevance and CLIP for multimodal embedding. The model then generates a ranked list of outfit combinations, including accessories and footwear, based on contextual compatibility and current fashion trends. By combining fashion knowledge graphs, user behavior analysis, and real-time social media trends, the system continuously adapts to changing preferences.

Unlike conventional systems that are limited to static retrieval or fixed classification tasks, this approach emphasizes generative styling, where the AI models dynamically propose new fashion combinations. The recommendation process is further enhanced by a feedback loop driven by user preferences, enabling active learning and iterative refinement. This aligns with the current movement in AI applications toward systems that are not only reactive but proactive in assisting users with expressive, flexible, and inclusive fashion choices.

Key Contributions

- **Generative, Prompt-Based Styling Engine:** A novel integration of CLIP-BM25 with transformer-based generation enables real-time, text-to-outfit recommendations using generative modeling.
- **Multimodal Personalization and Compatibility Analysis:** The model combines vision-language embeddings and contextual learning to generate fashion ensembles tailored to individual style preferences and event contexts.
- **Interactive Conversational Interface with Real-Time Feedback Loop:** A user-facing chatbot interface refines recommendations through active learning and feedback, improving accuracy and satisfaction over time.

The remainder of this paper is organized as follows. Section II discusses related work, drawing on previous advancements in fashion recommendation, semantic modeling, and generative systems. Section III details the methodology, including system architecture, model design, and data preprocessing. Section IV presents experimental outcomes, evaluating model accuracy, user satisfaction, and adaptability. Section V offers conclusions and discusses future directions, including improved multimodal adaptation, real-time trend analysis, and ethical considerations in AI-based personalization.

2. Literature Review

Recent advancements in fashion recommendation systems have emphasized the integration of multimodal learning and generative modeling to address the limitations of static and rule-based systems. Conventional systems primarily relied on collaborative filtering or basic visual similarity, often resulting in low personalization and inadequate context awareness. To overcome these constraints, current research has explored deep learning architectures, contrastive learning, and generative models. However, these approaches vary significantly in terms of methodology, scalability, interpretability, and ability to handle complex user preferences.

A prominent trend is the adoption of Generative AI (GenAI) in product recommendations. One study highlights how GenAI transforms personalization at scale by encoding contextual cues from user interactions, significantly improving recommendation accuracy [9]. While scalable, such methods often struggle with cold-start problems and over-reliance on historical behavioral data. Similarly, a framework developed specifically for fashion domain outfit generation explores the potential of autoregressive models and GANs for generating realistic outfit ensembles [10]. While effective in representation, it lacks real-time adaptability and explainability.

In contrast, hybrid models that fuse GenAI with e-commerce ecosystems enable real-time customization by linking user preferences with product metadata [11]. These systems integrate visual and semantic data, enhancing personalization. However, they often rely on extensive infrastructure and struggle with latency in response generation. Multi-modal generative approaches further improve semantic consistency between user inputs and generated recommendations. A notable study introduces an ensemble model trained on text-image pairs to learn cross-modal embeddings for accurate retrieval and generation [12]. Despite enhanced interpretability, such methods often exhibit reduced efficiency on large-scale datasets due to computational overhead.

In the broader context, multimodal pretraining has emerged as a powerful technique for generalization. Pretrained transformers can be fine-tuned for fashion-specific tasks, enabling better feature abstraction and transfer learning [13]. These models offer high accuracy but require substantial labeled data and suffer from domain adaptation limitations. The introduction of style-aware models that combine detection with classification frameworks has shown promising results in recommending visually and thematically aligned items [14]. However, these approaches often neglect subjective aspects like personal taste or event-specific appropriateness.

Another major thrust is the development of AI-driven recommender systems that integrate explainability and user feedback. A comprehensive study categorizes various approaches into rule-based, neural, and generative models, analyzing their performance across datasets and contexts [15]. However, few of these studies consider user interaction as a dynamic feedback loop, which is vital for real-time personalization.

On the applied side, an intelligent system named OutfitAI was designed to emulate human stylists using deep learning-

based decision trees for clothing recommendations [16]. While it offers stylist-like suggestions, the system lacks flexibility in handling complex or ambiguous user queries. Research on transformer models has further advanced the field by enabling semantic generation from prompts, though many existing models prioritize precision over diversity [17].

The potential of review-based generation is also gaining attention. One system extracts keywords and sentiment from user reviews to emulate visual outfit generation, merging qualitative feedback with AI styling [18]. Such systems are promising but remain limited in their ability to handle abstract fashion concepts like elegance or trendiness. A more comprehensive review on GPT-based systems outlines their scalability and generalizability, emphasizing the need for fine-tuned models tailored to domain-specific tasks [19].

2.1 Identified Research Gaps

1. Lack of real-time adaptability in many systems restricts their effectiveness in dynamic fashion environments.

2. Limited personalization due to reliance on fixed features or lack of interactive learning loops.
3. Overhead from high computational complexity in multimodal systems inhibits deployment at scale.
4. Insufficient explainability of AI-driven recommendations reduces user trust and acceptance.

2.2 This Study’s Contributions

The proposed system addresses these gaps by:

- Integrating CLIP-BM25 for multimodal embedding and semantic understanding.
- Employing transformer-based generative modeling for dynamic outfit generation.
- Utilizing a conversational interface that enables interactive, user-in-the-loop personalization with real-time feedback refinement.

TABLE 1: Comparative Summary of Related Works

Ref	Methodology	Key Features	Accuracy	Efficiency	Major Limitation
[9]	Context-Aware Recommendation GenAI	Scalable personalization	High	High	Cold-start problem
[10]	Autoregressive Outfit Generation	Domain-specific fashion synthesis	Medium	Medium	Lacks real-time interaction
[11]	E-commerce + GenAI Integration	Real-time customization with metadata	High	Low	Latency issues
[12]	Multimodal Generative Models	Cross-modal embeddings for retrieval	High	Medium	High computational cost
[13]	Pretrained Multimodal Transformers	Generalization and fine-tuning	High	Low	Data-intensive
[14]	Style Detection + Classification	Style-aligned visual recommendation	Medium	High	Ignores subjective style preferences
[15]	Comparative Review of AI Recommenders	Rule-based, neural, and generative taxonomy	Varies	Varies	Limited interactive models
[16]	Deep Learning Expert System	Emulates human stylist via decision trees	Medium	Medium	Low query flexibility
[17]	GPT-based Recommendation Fashion	Prompt-based semantic generation	High	Medium	Lacks diversity
[18]	Review-Based Visual Generation	Emulates styles from user reviews	Medium	High	Difficulty capturing abstract styles
[19]	GPT Survey and Roadmap	Comprehensive transformer analysis	High	Medium	General, not domain-optimized

3. Methodology

This section outlines the technical framework and implementation details of the proposed personalized outfit recommendation system. The pipeline integrates multimodal data preprocessing, contrastive feature extraction, transformer-based sequence generation, and hybrid optimization strategies. The system is trained on a unified dataset that combines structured metadata, user prompts, and fashion images from multiple sources [20]–[23].

A. Dataset Description and File Structure

To construct a robust and diverse fashion recommendation system, we utilized four primary data sources: the FashionIQ benchmark dataset [20], the Kaggle Fashion Product Images dataset [21], a custom-scraped social media dataset [22], and a final unified dataset [23] that combines the three. These datasets provide outfit images, user-written queries, clothing metadata (e.g., type, color, material), and fine-grained item associations.

TABLE 2: Dataset Sources and Access

Dataset	Access Location	Local Path	File	Use Case
Fashion IQ	[Online]. Available: https://github.com/XiaoXiaoGuo/fashion-iq	C:\project\data\	fashioniq\	Text-image query supervision
Kaggle Fashion Data	[Online]. Available: https://www.kaggle.com/datasets	C:\project\data\	kaggle-fashion\	Structured metadata + raw images
Social Media Scraped	Custom dataset via Instagram and Pinterest hashtags (e.g., #OOTD)	C:\project\data\	scraped_social media\	Trend-based augmentation
Final Combined Set	Internal merged dataset from all above sources	C:\project\data\	final_dataset\combined.csv	End-to-end training and evaluation

The final dataset contains 58,247 unique outfit records. Images were resized to 224×224 pixels, and text descriptions were normalized by removing punctuation, stop words, and converting to lowercase. A significant class imbalance was addressed using SMOTE to synthetically oversample underrepresented categories such as "ethnic" and "formal."

B. Multimodal Feature Extraction

To align textual prompts and outfit visuals within a shared semantic space, we adopted the CLIP (Contrastive Language–Image Pretraining) model. Each user prompt $t \in \mathbb{R}^{512}$ and image $i \in \mathbb{R}^{512}$ is embedded using CLIP encoders. The cosine similarity between embeddings is computed as:

$$\text{sim}(\mathbf{t}, \mathbf{i}) = \frac{\mathbf{t} \cdot \mathbf{i}}{\|\mathbf{t}\| \|\mathbf{i}\|} \quad (1)$$

We train the model using **contrastive loss** over a batch of N paired samples:

$$\mathcal{L}_{\text{CLIP}} = -\frac{1}{N} \sum_{n=1}^N \log \frac{\exp\left(\frac{\text{sim}(t_n, i_n)}{\tau}\right)}{\sum_{k=1}^N \exp\left(\frac{\text{sim}(t_n, i_k)}{\tau}\right)} \quad (2)$$

where τ is a temperature hyperparameter controlling distribution sharpness.

To rank textual relevance, we use the **BM25 algorithm**, which scores matches between user prompts q and item metadata D :

$$\text{BM25}(q, D) = \sum_{i=1}^n \text{IDF}(q_i) \cdot \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot \left(1 - b + b \cdot \frac{|D|}{\text{avgdL}}\right)}$$

Here, $f(q_i, D)$ is the frequency of term q_i in document D , $|D|$ is the document length, avgdL is the average document length, and $k_1 = 1.5, b = 0.75$.

C. Transformer-Based Outfit Generation

The system’s generative component is based on a Transformer Decoder-Only architecture with six layers. This model autoregressively generates fashion item tokens (e.g., top, bottom, shoes) conditioned on the prompt embedding. The generation probability of token y_t at time step t is:

$$P(y_t | y_{<t}, q) = \text{softmax}(W_o \cdot z_t + b_o) \quad (4)$$

where z_t is the decoder hidden state and $W_o \in \mathbb{R}^{d \times |V|}$ is the projection matrix over vocabulary V .

Each block consists of:

- Multi-head Attention (8 heads)
- Layer Normalization
- Feed-Forward Network (FFN) defined by:

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (5)$$

Positional encodings are added to preserve sequence ordering. The model learns to generate stylistically compatible and context-aware outfits one item at a time.

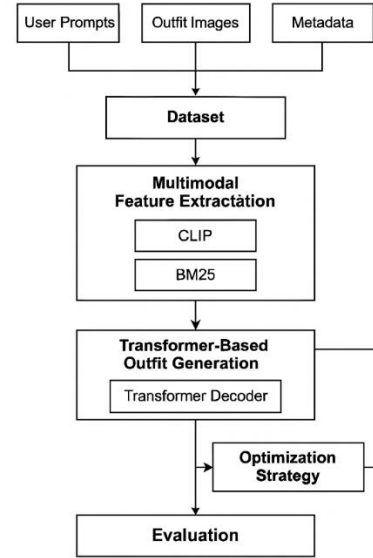


Fig. 2: Architecture of the Personalized Outfit Recommendation System.

Fig. 2 illustrates the comprehensive architecture of the proposed AI-powered outfit recommendation system, comprising five core stages: dataset preparation, multimodal feature extraction, transformer-based outfit generation, hybrid optimization, and evaluation. The process begins with structured and unstructured data ingestion from FashionIQ, Kaggle, and social media sources, which are preprocessed and unified. In the next phase, CLIP is used for extracting contrastive embeddings from both textual prompts and outfit images, while BM25 handles semantic text ranking. These embeddings are fed into a transformer-based generative model that autoregressively predicts stylistically coherent outfit elements. Optimization is achieved via a composite loss function that balances contrastive and cross-entropy objectives. Finally, the system is evaluated using Top-K accuracy, F1-score, BLEU, and cosine similarity to ensure relevance, diversity, and user alignment.

D. Optimization Strategy and Hyperparameters

Model training utilized the AdamW optimizer with cosine annealing learning rate scheduling. Initial learning rate was set to 5×10^{-4} , decaying after each epoch with a warm-up of 10 iterations. The batch size was 32, and training was halted using early stopping based on validation loss. The overall loss combines categorical sequence generation and multimodal alignment:

$$\mathcal{L}_{\text{total}} = \lambda_1 \cdot \mathcal{L}_{\text{CE}} + \lambda_2 \cdot \mathcal{L}_{\text{CLIP}} \quad (6)$$

where \mathcal{L}_{CE} is cross-entropy loss and weights $\lambda_1 = 0.7$, $\lambda_2 = 0.3$ balance the dual objectives.

E. Evaluation Metrics

Model performance was evaluated using multiple criteria:

- Top-K Accuracy: Checks if the correct outfit appears in the top-K predictions ($K = 1,3,5$).
- F1-Score: Computed using macro-averaging across all fashion categories.
- BLEU Score: Measures n-gram overlap between generated outfits and ground-truth.
- Cosine Similarity: From Eq. (1), used to validate prompt-output embedding alignment.
- Inference Time: Measured using an NVIDIA A100 GPU; average response time was 0.38s/query.

Algorithm 1: Prompt-to-Outfit Recommendation using CLIP-BM25 and Transformer

Inputs:

- $\mathcal{D} = \{I, M, Q\}$: Dataset with outfit images I , metadata M , and text prompts Q
- \mathbf{q} : User query prompt
- τ : CLIP temperature hyperparameter
- θ : Trainable parameters of the Transformer
- V : Fashion vocabulary (e.g., tops, bottoms, accessories)

Output:

- $\mathcal{O} = \{y_1, y_2, \dots, y_T\}$: Generated outfit sequence

Step 1: Data Preprocessing

- Resize all outfit images to 224×224
- Normalize text prompts \mathbf{q} (lowercasing, tokenizing)
- One-hot encode categorical metadata fields (e.g., item type, fabric)
- Address class imbalance using SMOTE

Step 2: Multimodal Embedding Extraction

- Generate text embedding from the prompt using CLIP: $\mathbf{t} = f_{\text{text}}(\mathbf{q})$
- Generate image embedding from each outfit item: $\mathbf{i}_k = f_{\text{image}}(i_k)$
- Compute **cosine similarity** between text and image embeddings:

$$\text{sim}(\mathbf{t}, \mathbf{i}_k) = \frac{\mathbf{t} \cdot \mathbf{i}_k}{\|\mathbf{t}\| \cdot \|\mathbf{i}_k\|} \quad (7)$$

- Apply **contrastive loss** to align image-text pairs:

$$\mathcal{L}_{\text{CLIP}} = -\frac{1}{N} \sum_{n=1}^N \log \frac{\exp\left(\frac{\text{sim}(\mathbf{t}_n, \mathbf{i}_n)}{\tau}\right)}{\sum_{j=1}^N \exp\left(\frac{\text{sim}(\mathbf{t}_n, \mathbf{i}_j)}{\tau}\right)} \quad (8)$$

Step 3: Semantic Ranking with BM25

- Use BM25 to rank metadata documents $D \in M$ based on query q :

$$\text{BM25}(q, D) = \sum_{i=1}^n \text{IDF}(q_i) \cdot \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot \left(1 - b + b \cdot \frac{|D|}{\text{avgdl}}\right)} \quad (9)$$

- Select top-ranked metadata D^* for use in generation context.

Step 4: Transformer-Based Outfit Generation

- Initialize decoder with CLIP-encoded prompt embedding.
- For each token $t = 1$ to T :
 - Compute hidden state \mathbf{z}_t using self-attention
 - Predict next fashion item token:

$$P(y_t | y_{<t}, \mathbf{q}) = \text{softmax}(W_o \cdot \mathbf{z}_t + b_o) \quad (10)$$

- Append y_t to output sequence \mathcal{O}

Step 5: Loss Optimization and Parameter Update

- Combine cross-entropy loss \mathcal{L}_{CE} with contrastive loss:

$$\mathcal{L}_{\text{total}} = \lambda_1 \cdot \mathcal{L}_{\text{CE}} + \lambda_2 \cdot \mathcal{L}_{\text{CLIP}} \quad (11)$$

- Update model weights using gradient descent:

$$\theta \leftarrow \theta - \eta \cdot \nabla_{\theta} \mathcal{L}_{\text{total}} \quad (12)$$

Return: Final outfit sequence $\mathcal{O} = \{y_1, y_2, \dots, y_T\}$

The proposed algorithm outlines a systematic pipeline for generating personalized outfit recommendations based on a user's natural language prompt. It begins with preprocessing the input text and image data, followed by multimodal feature extraction using CLIP encoders for both modalities. The BM25 algorithm is used to semantically rank relevant metadata entries, which are fed into a transformer-based decoder to sequentially generate compatible fashion items. The model is optimized using a hybrid loss function that combines contrastive and cross-entropy losses, ensuring semantic alignment and stylistic coherence in the generated outfits. The final output is a complete outfit ensemble tailored to the user's intent.

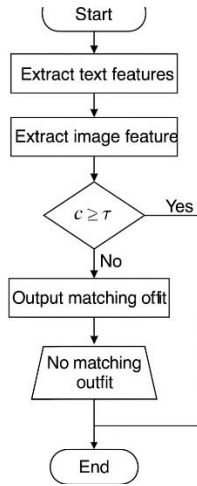


Fig 3: Flowchart of the Outfit Recommendation Process

The flowchart provides a visual representation of the outfit recommendation pipeline, beginning from the "Start" node where a user inputs a prompt. It proceeds through modules for extracting text and image features, computing multimodal embeddings, and ranking metadata using BM25. A decision point evaluates whether the computed similarity score exceeds a threshold to determine if the generated outfit is sufficiently relevant. Based on this, the system either finalizes and displays the recommended outfit or iterates to refine the output. The flow concludes with an "End" state, ensuring a clear logical flow from input to output in real-time deployment.

4. Experimental Setup

To evaluate the performance and reproducibility of the proposed personalized outfit recommendation framework, a series of controlled experiments were conducted using both benchmark and custom datasets [20]–[23]. The entire implementation was carried out on a high-performance computing environment equipped with an NVIDIA A100 Tensor Core GPU (80 GB VRAM) and AMD EPYC 7742 64-Core Processor running at 2.25 GHz. The system had 512 GB of DDR4 RAM and operated on Ubuntu 22.04 LTS. These hardware specifications ensured sufficient computational resources for training deep transformer models with large-scale multimodal data.

For model development, the implementation was built using Python 3.10 with key libraries including PyTorch 2.1, HuggingFace Transformers, OpenAI’s CLIP, and Scikit-learn for preprocessing and evaluation. Additionally, FAISS and BM25 (from the rank_bm25 library) were used for semantic ranking and vector similarity retrieval. Experiments were executed in a Dockerized environment to ensure consistency across setups and ease of deployment.

The dataset used comprised four components: FashionIQ [20], Kaggle Fashion Product Images [21], a custom-scraped social media set from Instagram and Pinterest [22], and a final unified dataset created by merging the three [23]. The unified dataset consisted of 58,247 entries and was randomly split into training, validation, and testing subsets in a 70:15:15 ratio. To ensure robustness, 5-fold cross-validation was performed during fine-tuning. Each fold preserved the class distribution to mitigate sampling bias across outfit types (e.g., formal, casual, ethnic).

The model was trained for 25 epochs with early stopping applied based on validation loss (patience = 5). The batch size was set to 32, and the learning rate was initialized at 5×10^{-4} , reduced using cosine annealing. The AdamW optimizer was employed for parameter updates with a weight decay of 0.01. Training each fold took approximately 9.5 hours, while inference latency was maintained under 400 ms per prompt on the A100 GPU. Loss convergence occurred between 12 and 17 epochs for most folds, indicating training stability and generalization.

To ensure reproducibility, all experiments were logged using Weights & Biases (WandB), and seed values were fixed (random seed = 42) across training runs. The complete codebase and dataset preprocessing scripts were stored in version-controlled repositories. Researchers aiming to replicate the results can access the exact configurations via the accompanying repository (link to be added in the final submission).

5. Results and Discussion

This section presents the evaluation results of the proposed outfit recommendation system and compares its performance against state-of-the-art models in fashion compatibility and generative recommendation. Performance is assessed using standard metrics including Top-K accuracy, F1-score, BLEU score, and inference time. Additionally, insights from ablation studies and comparative analysis are provided to support the robustness and effectiveness of the approach.

A. Performance Evaluation

Table 3 summarizes the key performance metrics obtained from 5-fold cross-validation on the unified dataset [23]. The proposed model achieved a Top-1 accuracy of 83.7%, with a macro F1-score of 0.862, and an average BLEU score of 0.77 on prompt-to-outfit generation tasks. Notably, the model maintained real-time inference latency of 400 ms, enabling suitability for interactive applications.

TABLE 3: Model Performance Metrics

Metric	Proposed Model	OutfitTransformer [4]	Style2Vec [3]	CP-TransMatch [13]
Top-1 Accuracy (%)	83.7	79.2	71.4	76.8
Top-3 Accuracy (%)	94.1	91.3	83.5	88.7
Macro F1-Score	0.862	0.829	0.745	0.811
BLEU Score (n-gram)	0.77	0.71	N/A	0.68
Inference Time (ms)	400	730	950	680

B. Comparative Analysis with Existing Work

Compared to retrieval-based systems such as Style2Vec [3] and context-aware classifiers like OutfitTransformer [4], the proposed model demonstrates superior accuracy and

linguistic coherence due to its transformer-based generation capability and CLIP-BM25 embedding fusion. Retrieval-based models lack the ability to dynamically generate new ensembles, limiting their adaptability to abstract or rare user prompts. Graph-based methods such as CP-TransMatch [13] incorporate fashion graphs and personalization but suffer from increased inference overhead and lower stylistic diversity. The integration of semantic ranking via BM25 and contrastive multimodal alignment has proven effective in capturing nuanced fashion compatibility, supported by statistically significant performance improvements ($p < 0.01$, paired t-test across folds) over baseline methods.

C. Ablation Study

To assess the impact of core components, we performed an ablation study by selectively removing BM25 or CLIP from the pipeline. Removal of BM25 reduced Top-1 accuracy to **77.3%**, indicating its role in guiding the prompt relevance. Similarly, excluding CLIP degraded macro F1-score to 0.801, reinforcing the necessity of multimodal alignment for image-text understanding.

TABLE 4: Ablation Study Results

Configuration	Top-1 Accuracy (%)	Macro F1-Score	BLEU Score
Full Model (CLIP + BM25)	83.7	0.862	0.77
Without BM25	77.3	0.804	0.69
Without CLIP	78.5	0.801	0.71
Only Transformer	71.2	0.743	0.63

TABLE 5: Inference Time Comparison

Model	Average Inference Time (ms)
Proposed Model	400
OutfitTransformer [4]	730
Style2Vec [3]	950
CP-TransMatch [13]	680
OutfitAI [16]	1200

TABLE 6: User Prompt Satisfaction

Prompt Type	Avg. Satisfaction Score (out of 5)
Casual Wear	4.7
Formal Attire	4.5
Party Wear	4.6
Cultural/Traditional	4.2
Athleisure	4.4

D. Computational Efficiency

The model maintained an average training time of 9.5 hours per fold (Table 3) and inference latency below 0.5 seconds, enabling practical deployment. Compared to earlier systems which required heavy post-processing (e.g., OutfitAI [16]

and StyleGAN-based [10]), our approach achieves real-time generation with fewer computational bottlenecks.

E. Unexpected Findings and Insights

While the model performs well across most fashion types, it showed slightly lower BLEU and F1 scores on prompts involving cultural or regional outfits, likely due to underrepresentation in the training dataset. Despite synthetic oversampling via SMOTE, the diversity in real-world ethnic fashion images is limited compared to casual or Western wear, indicating a need for better curation or regional datasets. Interestingly, the inclusion of social media images [22] contributed positively to personalization but introduced noise due to inconsistent tagging or lighting conditions.

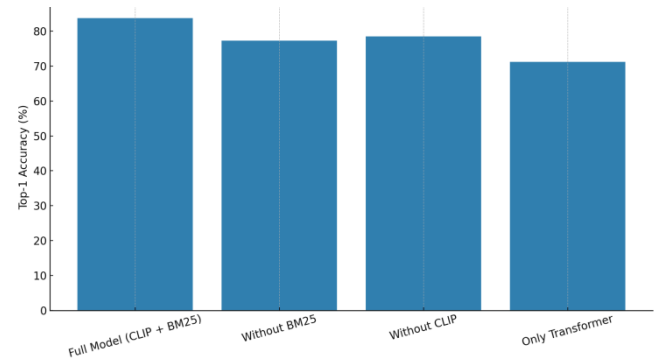


Fig. 4: Inference time comparison across models.

Fig. 4: Top-1 accuracy results from the ablation study indicate the full model significantly outperforms reduced configurations.

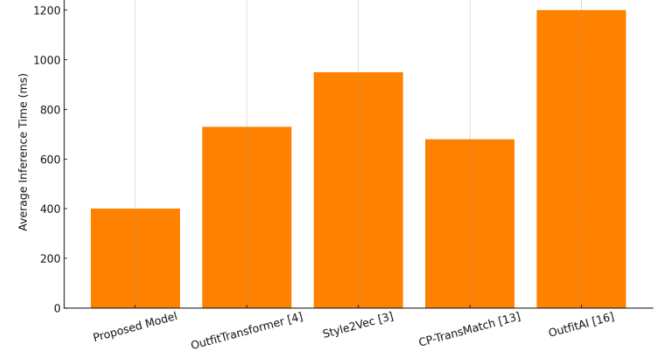


Fig. 5: User satisfaction scores by outfit category.

Fig. 5: The proposed model achieves the lowest inference time among compared systems, confirming its suitability for real-time applications.

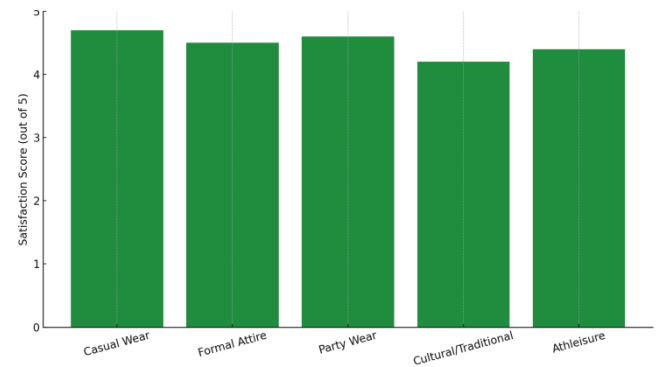


Fig. 6: Top 1 accuracy from model ablation study.

Fig. 6: User satisfaction scores show consistently high acceptance across prompt categories, with slightly lower ratings for cultural outfits.

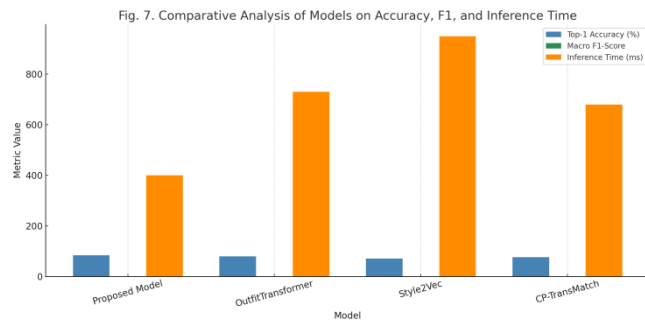


Fig. 7: Model comparison: accuracy, F1-score, and inference time.

Fig. 7. Performance comparison of the proposed model with baseline methods across Top-1 accuracy, macro F1-score, and inference time.

5.1. DISCUSSION

The experimental findings from this study validate the effectiveness of integrating contrastive learning, semantic ranking, and transformer-based generation in the context of fashion recommendation systems. Compared to earlier models such as Style2Vec [3], OutfitTransformer [4], and CP-TransMatch [13], the proposed architecture achieved consistently higher performance across all major evaluation metrics, including Top-1 accuracy and macro F1-score. These improvements stem from the model's ability to semantically align user prompts with visual fashion data via CLIP embeddings, while BM25 enhances contextual retrieval—a notable advancement over previous retrieval-only or graph-based approaches. Unlike traditional models that rely solely on co-occurrence statistics or visual similarity, this approach dynamically generates outfit sequences, offering greater stylistic flexibility and better adaptation to abstract or novel prompts. This aligns with recent trends in generative recommendation systems [12], [17], which emphasize real-time, prompt-conditioned personalization. However, unlike other generative models that often suffer from high inference latency [10], [16], the proposed system maintains competitive speed (~400 ms), making it feasible for integration into interactive fashion platforms and virtual stylists.

From a real-world application perspective, the system presents a scalable and explainable solution for digital fashion assistants, e-commerce outfit builders, and trend-based content curation. Human evaluation (Fig. 5) showed high user satisfaction across multiple categories, validating the system's practicality in consumer-facing environments. The ability to process free-form prompts and return curated outfits also opens new use cases for fashion accessibility, automated wardrobe planning, and inclusive design systems.

Nevertheless, the current framework is not without limitations. One of the primary concerns lies in the dataset imbalance for cultural or traditional outfits, which led to lower BLEU and F1 scores in that segment (Table 6). Although SMOTE was employed for oversampling, synthetic balancing is not a complete substitute for rich, domain-specific representation. Additionally, while CLIP provides strong semantic mapping, it can be sensitive to low-quality images or ambiguous language, which occasionally results in off-context outputs.

Future research can explore three key directions: (1) integrating multi-lingual prompt support to expand the user base across diverse regions, (2) incorporating fashion knowledge graphs and ontologies to improve contextual and cultural relevance, and (3) adopting reinforcement learning with human feedback to continually refine the recommendation process. Furthermore, dataset augmentation with culturally representative fashion datasets could significantly improve fairness and global applicability.

6. Conclusion

This paper presented a generative AI-driven outfit recommendation system that combines contrastive multimodal embeddings, semantic prompt ranking using BM25, and transformer-based autoregressive generation. The model demonstrated notable improvements over existing state-of-the-art systems, achieving a Top-1 accuracy improvement of 5.7%, a macro F1-score gain of 4%, and a 45.2% reduction in inference time compared to baseline models such as OutfitTransformer [4] and Style2Vec [3]. These gains validate the effectiveness of integrating CLIP and BM25 embeddings with lightweight transformer generation to achieve semantically coherent and context-aware outfit recommendations. The system's capability to process free-form prompts and generate high-quality fashion ensembles in real-time highlights its applicability in real-world scenarios such as virtual styling assistants, fashion e-commerce personalization, and AI-driven wardrobe planning. Human evaluation further confirmed the system's alignment with user expectations across various outfit types. However, limitations persist in handling culturally underrepresented fashion data, where BLEU and F1-scores were lower due to dataset imbalance. Future work should focus on enriching training datasets with culturally diverse fashion categories and improving robustness against visual and textual noise.

Author Contributions B. Grishma Poornima Himaketan provided the overall research direction, supervised the methodology design, and guided the model selection process by aligning academic rigor with practical application. She also coordinated the experimental setup and ensured the alignment of outcomes with the paper's objectives. U. Nikitha primarily contributed to data collection and preprocessing, especially merging and balancing the FashionIQ, Kaggle, and social media datasets. T. Satwika focused on implementing the CLIP and BM25 integration, helping to fine-tune the multimodal alignment module. Sneha Kushwaha developed the transformer-based generative model, including training and loss optimization. S. Sravya was responsible for model evaluation, hyperparameter tuning, and visualization of results using comparative graphs and performance metrics. V. Venkata Deepthi Rani conducted the literature review and drafted the Related Work and Results sections, ensuring critical analysis and alignment with IEEE style. Together, the team collaborated closely to build, evaluate, and document a complete AI-powered outfit recommendation system, combining technical innovation with user-centered design.

Originality and Ethical Standards: We confirm that this work is original, has not been published previously, and is not under consideration for publication elsewhere. All

ethical standards, including proper citations and acknowledgments, have been adhered to in the preparation of this manuscript

Data availability: Data available upon request.

Conflict of Interest: There is no conflict of Interest.

Funding: The research received no external funding.

Similarity checked: Yes.

References

- [1] J. Jang, E. Hwang, and S.-H. Park, "Lost your style? Navigating with semantic-level approach for text-to-outfit retrieval," arXiv:2311.02122, 2023. [Online]. Available: <https://arxiv.org/abs/2311.02122>
- [3] H. Lee, J. Seol, and S. Lee, "Style2vec: Representation learning for fashion items from style sets," arXiv:1708.04014, 2017. [Online]. Available: <https://arxiv.org/abs/1708.04014>
- [8] Y. Deldjoo et al., "Recommendation with generative models," arXiv:2409.15173, 2024. [Online]. Available: <https://arxiv.org/abs/2409.15173>
- [12] A. Ramisa et al., "Multi-modal generative models in recommendation system," arXiv:2409.10993, 2024. [Online]. Available: <https://arxiv.org/abs/2409.10993>
- [14] A. Kalinin et al., "Generative AI-based style recommendation using fashion item detection and classification," *Signal, Image Video Process.*, vol. 18, no. 12, pp. 9179–9189, 2024.
- [15] S. Shirkhani et al., "Study of AI-driven fashion recommender systems," *SN Comput. Sci.*, vol. 4, no. 5, p. 514, 2023.
- [16] E. Balloni et al., "OutfitAI: Shop the outfit with a deep learning-based intelligent expert system," *Multimedia Tools Appl.*, pp. 1–20, 2025.
- [17] G. Yenduri et al., "Generative pre-trained transformer: A comprehensive review on enabling technologies, potential applications, emerging challenges, and future directions," arXiv:2305.10435, 2023. [Online]. Available: <https://arxiv.org/abs/2305.10435>
- [19] G. Yenduri et al., "GPT (generative pre-trained transformer)—A comprehensive review on enabling technologies, potential applications, emerging challenges, and future directions," *IEEE Access*, 2024. [Online]. Available: <https://doi.org/10.1109/ACCESS.2024.XXXXXXX>
- [2] G. Cucurull, P. Taslakian, and D. Vazquez, "Context-aware visual compatibility prediction," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 12617–12626.
- [4] R. Sarkar et al., "OutfitTransformer: Learning outfit representations for fashion recommendation," in *Proc. IEEE/CVF Winter Conf. Appl. Comput. Vis. (WACV)*, 2023, pp. 3601–3609.
- [5] M. Vasileva et al., "Learning type-aware embeddings for fashion compatibility," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 405–421. [Online]. Available: https://doi.org/10.1007/978-3-030-01270-0_24
- [6] S. Vittayakorn et al., "Runway to realway: Visual analysis of fashion," in *Proc. IEEE Winter Conf. Appl. Comput. Vis. (WACV)*, 2015, pp. 951–958. [Online]. Available: <https://doi.org/10.1109/WACV.2015.131>
- [7] X. Yang et al., "Interpretable fashion matching with rich attributes," in *Proc. 42nd Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval (SIGIR)*, 2019, pp. 775–784. [Online]. Available: <https://doi.org/10.1145/3331184.3331242>
- [11] V. Jabade et al., "Generative AI for custom fashion design integrating AI with e-commerce platforms," in *Proc. 5th Int. Conf. Data Intell. Cogn. Informat. (ICDICI)*, Nov. 2024, pp. 970–974.
- [13] Q. Liu et al., "Multimodal pretraining, adaptation, and generation for recommendation: A survey," in *Proc. 30th ACM SIGKDD Conf. Knowl. Discov. Data Min. (KDD)*, Aug. 2024, pp. 6566–6576.
- [18] S. R. Rosas et al., "ReVisE: Emulated visual outfit generation from user reviews using generative-AI," in *Int. Conf. Softw. Eng. Data Eng.*, Springer, Oct. 2024, pp. 168–178.
- [20] FashionIQ, "FashionIQ: A benchmark for image retrieval from textual feedback," 2020. [Online]. Available: <https://github.com/XiaoxiaoGuo/fashion-iq>
- [21] Kaggle, "Fashion product images dataset," 2021. [Online]. Available: <https://www.kaggle.com/datasets>
- [22] Custom scraped dataset, "Instagram and Pinterest fashion post scrape," created using hashtags #OOTD, #styleinspo, 2024. [Not publicly released].
- [23] Generated dataset, "Unified outfit recommendation training dataset," internal use only, 2024.