# FashionDPO: Fine-tune Fashion Outfit Generation Model using Direct Preference Optimization

Mingzhe Yu
mz_y@mail.sdu.edu.cn
Shandong University
Jinan, China

Yunshan Ma
ysma@smu.edu.sg
Singapore Management University
Singapore

Lei Wu*
i_lily@sdu.edu.cn
Shandong University
Jinan, China

Changshuo Wang
cswang02@outlook.com
Shanghai Jiao Tong University
Shanghai, China

Xue Li
xue.lii754@gmail.com
Jiangnan University
Wuxi, China

Lei Meng
lmeng@sdu.edu.cn
Shandong University
Jinan, China

## Abstract

Personalized outfit generation aims to construct a set of compatible and personalized fashion items as an outfit. Recently, generative AI models have received widespread attention, as they can generate fashion items for users to complete an incomplete outfit or create a complete outfit. However, they have limitations in terms of lacking diversity and relying on the supervised learning paradigm. Recognizing this gap, we propose a novel framework FashionDPO, which fine-tunes the fashion outfit generation model using direct preference optimization. This framework aims to provide a general fine-tuning approach to fashion generative models, refining a pre-trained fashion outfit generation model using automatically generated feedback, without the need to design a task-specific reward function. To make sure that the feedback is comprehensive and objective, we design a multi-expert feedback generation module which covers three evaluation perspectives, *i.e.,* quality, compatibility and personalization. Experiments on two established datasets, *i.e.,* iFashion and Polyvore-U, demonstrate the effectiveness of our framework in enhancing the model's ability to align with users' personalized preferences while adhering to fashion compatibility principles. Our code and model checkpoints are available at https://github.com/Yzcreator/FashionDPO.

## CCS Concepts

• **Information systems → Multimedia and multimodal retrieval**; **Recommender systems**; **Multimedia content creation**.

## Keywords

Fashion Outfit Generation, Fashion Image Generation, Generative Fashion Recommendation

## 1 Introduction

Fashion in life is not merely about clothing and dressing, while it serves as a means of expressing personal style, identity, and lifestyle [2]. One of the major requirements in people's daily fashion life is personalized fashion matching, where people need to pick up a set of fashion items (defined as an outfit) that are compatible with each other and fit his/her personalized fashion taste [43, 48]. In addition to experts in the domain of fashion, researchers in information retrieval also express interest in studying this problem, which is defined as the task of personalized outfit recommendation [1, 5, 24]. The specific tasks include 1) Personalized Fill-in-the-Blank (PFITB), completing an incomplete outfit by finding a compatible fashion item based on the user's interaction history, and 2) Generative Outfit Recommendation (GOR), assembling a complete outfit from scratch catering to both personalized preference and fashion item compatibility.

Regarding these specific tasks, existing works can be categorized into retrieval-based methods and generation-based methods. The retrieval-based models [1, 5, 9, 50] typically explore the existing candidate item set and select certain items to either complete a partial outfit or curate an outfit from scratch. Despite their prevailing usage, such retrieval-based methods are inevitably restricted by the limited inclusivity and diversity of existing items. [35, 48] What if we cannot find a suitable item after iterating every item within the available item set? To address this problem, a natural and reasonable solution is to directly generate an image of the desired fashion item using the cutting-edge image generation models [15, 28, 34, 43, 51]. For example, DiFashion [43] leverages the pre-trained stable diffusion model and incorporates users' historical preference and items' compatible constraints, achieving SOTA performance on both tasks of PFITB and GOR. Despite its promising performance, we pinpoint a problem that the generated fashion items exhibit limited diversity, which is probably due to the fact that the pre-defined outfits in

**Figure 1: Illustration of our motivation and the paradigm comparison between FashionDPO and the supervised learning methods, where FashionDPO optimizes the model using feedback from multiple experts without relying on labeled dataset, resulting in high diversity while maintaining high quality.**

the dataset cannot fully cover all possible fashion item combinations. For example, in the PFITB task, each incomplete outfit is paired with only one ground-truth item, which, however, may have multiple alternative items that can also satisfy the personalization and compatibility constraints. Furthermore, current generation-based methods solely rely on the supervised learning paradigm, which exacerbates the issue of lacking diversity during the generation [40, 45]. As shown in Figure 1, during the training process, the generative models [28, 43] aim to minimize the difference between the result image and the ground-truth image. Thereby, when the generated images visually adhere to fashion compatibility but differ from the ground truth, the model would erroneously penalize these deviations. Consequently, the model would overly focus on specific features within the pre-defined outfits, rather than learning more generalized principles of clothing compatibility, resulting in limited generation diversity.

To resolve the limitations of lacking diversity and supervised learning paradigm, intuitive solutions could be augmenting the dataset through additional human annotations. However, relying on human annotations yield high monetary cost and time expense, even worse, the involvement of human annotators would introduce uncontrollable subjective biases into the dataset. Additionally, it becomes challenging to adapt or update the dataset when the fashion trends are constantly changing and evolving [4, 26]. Considering the above concerns, a more advanced and promising approach is the RLAIF [19] (reinforcement learning from AI feedback), which enables the model to learn from a variety of possible feedback and adapt to more diverse situations. This method is particularly suitable for scenarios, such as artistic creation, where no fixed answers exist. Nevertheless, how to design the AI models to ensure that the feedback is both comprehensive and objective still remains a challenge. Moreover, how to leverage the feedback generated by AI models to guide the training of generation models is still challenging and under-explored in the fashion domain. For example, existing methods [18] typically first train a reward model that is aligned with human preferences and then fine-tune the generation model using reinforcement learning. However, constructing an effective reward model demands large-scale datasets and multiple rounds of parameter tuning, which is time-consuming and resource-intensive.

To tackle these challenges, we propose a novel framework FashionDPO, which fine-tunes the fashion outfit generation model using direct preference optimization. To make sure the automatic feedbacks are comprehensive and professional, we design a multi-expert feedback generation module to address the first challenge, which consists of three experts covering diverse evaluation perspectives: 1) Quality: evaluating whether the fashion elements in the image are complete and conform to fashion design principles, 2) Compatibility: assessing whether the generated fashion products are coordinated within an outfit, and 3) Personalization: ensuring that the recommended fashion products align with the user's personal preferences. To solve the second challenge, we leverage the DPO (Direct Preference Optimization) [32] framework and group the feedbacks into positive-negative pairs, which are collected from multiple experts' feedback. We then utilize the positive-negative pairs to guide the fine-tuning process, eliminating the need to train a reward model. In summary, **the primary contributions** of our work are as follows:

- We identify the limitations of lacking diversity and supervised learning paradigm in existing personalized outfit generation models, and we propose to utilize mulitple AI feedbacks to solve these limitations.
- Fulfilling the idea of using AI feedbacks, we propose a novel framework FashionDPO, which fine-tunes the generation model using DPO based on multiple AI experts' feedbacks.
- Extensive experiments on iFashion and Polyvore-U datasets demonstrate the effectiveness of our framework regarding various evaluating metrics.

## 2 Related Work

**Fashion Outfit Recommendation.** In the fashion domain, Outfit Recommendation (OR) [20, 24, 50] has gained widespread application. There are two requirements in fashion outfit recommendation: compatibility and personalization. Furthermore, it is also a popular task in the domain of computational fashion [2]. Early works [6, 24, 50] primarily focused on compatibility, aiming to retrieve already well-matched outfits for users. Some works [1, 5] attempt to introduce personalization in the recommendation process, combining a set of personalized and compatible items into an outfit that aligns with fashion styling principles. Moreover, bundle recommendation, a more generalized recommendation paradigm, subsume personalized fashion outfit recommendation as one of its applications. Multiple works [7, 23, 27] have been proposed by using graph learning, contrastive learning, as well as multimodal large language models. Despite various progress, the above works follow the retrieval paradigm and are constrained by the variety and quantity of fashion products in the dataset, making it difficult to meet users' personalized needs, especially in terms of texture and other details. However, with the rapid development of generative models [34, 44, 51], the quality and diversity of image generation have significantly improved, making it possible to directly recommend generated custom fashion products to users. Recent work [43] has introduced the PFITB task, which combines the user's interaction history with fashion products to generate a personalized matching outfit.

**Fashion Image Generation.** It refers to the task of generating fashion-related images using deep learning models. This task is widely applied in the fashion domain, covering areas such as clothing design, virtual try-on, and personalized recommendation, among others [3, 35, 46]. Previous works, such as CRAFT [15], generate feature representations for clothes pairings and retrieve the most suitable individual clothes items from the dataset. In the virtual try-on domain, previous works [10, 42] based on GANs involve generating warped clothes aligned with character, and then generating images of the character wearing the warped clothes. The diffusion models [8] enhance image quality by replacing the generator in the second stage. Current work [16] learns the semantic correspondence between the clothing and the human body within the latent space of the pre-trained diffusion model in an end-to-end manner. In the personalized recommendation domain, HMaVTON [48] generates diverse and well-matched fashion items to the given person. Existing personalized image generation models [14, 29, 41, 49] aim to generate images aligned with reference styles or elements, yet recommending images consistent with a user's interaction history is meaningless.

**Direct Preference Optimization.** In the field of natural language processing, Direct Preference Optimization (DPO) has been proposed to reduce training costs [32], which uses preferences rather than explicit rewards to fine-tune LLMs. This approach is also applied to the post-training of text-to-image diffusion models. Diffusion-DPO [40] fine-tunes the generative model in a single step after receiving feedback from the Preference Evaluator. D3PO [45] assumes that the preferred outcome holds true for all time steps in the diffusion model and fine-tunes each of the time steps in the generative model based on the feedback results. It demonstrates that in diffusion models, directly updating the policy based on human preferences within an MDP is equivalent to first learning the optimal reward model and then using it to guide policy updates. SPO [21] assesses preferences at each time step during the sampling process and adds noise to the preferred image to generate the noise image for the next time step. We introduce DPO into generative fashion recommendation, where learning based on preference feedback eliminates the constraints of ground truth, showcasing richer possibilities in clothing textures and details.

## 3 Preliminary

We first introduce the problem formulation, followed by a briefing of the diffusion model and the direct preference optimization.

**Problem Formulation of PFITB and GOR.** Based on the user information $u$, these tasks aim to generate complete outfit $O = \{i_k\}_{k=1}^n$, where $O$ denotes complete outfit and $i_k$ is individual fashion item. Specifically, the PFITB task generates compatible fashion items $i_k$ for each incomplete outfit $O' = O \setminus \{i_k\}$ according to user preferences. And the GOR task is further expanded to generate a complete set of matching outfits $O$ for users.

**Diffusion Models.** Diffusion Models learn probability distribution $p(i)$ by inverting a Markovian forward process $q(i_t|i_{t-1})$. The GOR task, starting from multiple randomly initialized noisy images $O_T = \{i_{k,T}\}_{k=1}^n$ where $i_{k,T} \sim \mathcal{N}(0, I)$, the diffusion models gradually remove noise from the image through a reverse process. At time step $t \in \{1, ..., T\}$, the U-Net predicts the noise $\epsilon_k$ of the $i_{k,t}$, ultimately

generating personalized fashion outfits $O_0 = \{i_{k,0}\}_{k=1}^n$. And the PFITB task starts with a single noisy image $i_{k,T}$ and generates a fashion item $i_{k,0}$ for incomplete outfit $O'$. To achieve personalized outfit generation, we introduce the user's interaction history $h_k$, category prompt $p_k$ and mutual influence $m_k$ between items in the outfit as condition $c_k$. In the case of conditional generative modeling, the U-Net learns the probability distribution $p(i|c_k)$ with the following objective:

$$\mathcal{L} = \mathbb{E}_{t, \epsilon_k \sim \mathcal{N}(0,I)} [||\epsilon_k - \epsilon_\theta(i_{k,t}, h_k, p_k, m_k, t)||_2^2]. \tag{1}$$

**Direct Preference Optimization.** The purpose of preference-Based Models is to optimize the model's output according to the user's preferences. It often requires learning a reward function $r : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ from human feedback, where $\mathcal{S}, \mathcal{A}$ denotes state space and action space. The objective is to find a policy $\pi(a|s)$ that maximizes the reward $r$, where $a \sim \mathcal{A}, s \sim \mathcal{S}$. However, user preferences are typically complex and multi-dimensional, making it a challenge to accurately model and design the corresponding reward function. Recently, the direct preference optimization (DPO) [32] has been proposed to fine-tune models using direct feedback from users. Given a state $s$ and action pairs $(a_w, a_l)$, the $\pi_\theta(a|s)$ represents the sampling policy of the finetuned model, and $\pi_{ref}(a|s)$ represents the policy of the reference model. In the GOR task, we can generate fashion pairs $(y_w, y_l)$ based on prompt $x$, where $y_w$ is user's preference choice and $y_l$ is the inferior one. LPO formulates a maximum likelihood objective for a parametrized policy $\pi_\theta$:

$$\mathcal{L} = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta_1 \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta_2 \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right], \tag{2}$$

where $\beta$ is used to control the deviation between $\pi_\theta$ and $\pi_{ref}$.

## 4 Method

We present our proposed method FashionDPO, which consists of three consecutive modules: 1) fashion image generation without feedback, 2) feedback generation from multiple experts, and 3) model fine-tuning with direct preference optimization. We explicate the detailed design of these modules one by one.

### 4.1 Fashion Image Generation without Feedback

We faithfully follow DiFashion [43] to build the generation model. Starting from empty image $i_{1,0}$ incomplete outfit $O' = \{i_{k,0}\}_{k=2}^n$, where the $i_{k,0}$ represents a fashion item in the outfit. During the initialization process, the model adds random Gaussian noise to image $i_{1,0}$. The forward formula is expressed as follows:

$$i_{k,t} = \sqrt{\alpha_t} i_{k,0} + \sqrt{1 - \alpha_t} \epsilon_k, \tag{3}$$

where $\epsilon_k \sim \mathcal{N}(0, 1)$ is Gaussian noise sampled from a standard normal distribution, and $\alpha_t$ is a decay coefficient with time step $t \in \{1, 2, ..., T\}$. As $t$ increases, the influence of noise $\epsilon$ gradually grows, and $i_{k,T}$ becomes almost the entire noise. To obtain multiple generated results for extracting preference information, we apply the forward noising process $m$ times to obtain $I_T = \{i_{k,j,T}\}_{j=1}^m$.

In the reverse process, the model starts with noisy images $i_{k,j,t}$ as input, gradually predicts and removes the noise step by step, and ultimately reconstructs the fashion items $i_{k,j,0}$. The denoising process is outlined as follows:

**Figure 2: The overview of FashionDPO, which consists of three consecutive key modules: 1) Fashion Image Generation without Feedback, 2) Feedback Generation from Multiple Experts, and 3) Model Fine-tuning with Direct Preference Optimization.**

$$i_{k,j,0} = \frac{1}{\sqrt{\alpha_t}} (i_{k,j,t} - \sqrt{1 - \alpha_t} \epsilon_\theta(i_{k,j,t}, t, c_k)), \qquad (4)$$

where the $\epsilon_\theta(i_{k,j,t}, t, c_k)$ represents the noise predicted by the model, and $c_k$ denotes the condition. At each time step $t$, the model first saves the current latent images for the subsequent fine-tuning stage. As shown in Figure 2, the model introduces the mutual influence $m_k$, user's interaction history $h_k$ and category prompt $p_k$ as condition $c_k$:

**Mutual Influence.** To ensure that the generated fashion items within an outfit are well-matched, the model introduces mutual influence $m_k$ during the inference process, enabling the fashion items to influence each other. At time step $t$, the model uses the incomplete outfit $O' = \{i_{k,0}\}_{k=2}^{4}$ to obtain the latent feature representation $m_{1,t}$ through the Multi-Layer Perceptron (MLP), providing outfit matching information for $i_{1,j,t}$. And it fuses two latent representations $m_{1,t}$ and $i_{1,j,t}$ through element-wise addition:

$$i_{1,j,t}^m = (1 - \eta) \cdot i_{1,j,t} + \eta \cdot m_{1,t}, \qquad (5)$$

where $\eta$ is a hyperparameter that controls the influence of mutual condition. A larger $\eta$ indicates a greater impact on the noisy image.

**User's Interaction History.** The personalized information of the user is reflected through their fashion items' interaction history. For user $u$, the fashion items in the interaction history that belong to the same category as $i_{1,0}$ are extracted to form $u_k = \{i_k^h\}_{h=1}^N$. And the model use the pre-trained autoencoder $\mathcal{E}$ in diffusion models [34] to encode the user's interaction history images $u_k$ and take the average as the history condition $h_k$:

$$h_k = \frac{1}{N} \sum_{h=1}^{N} \mathcal{E}(i_k^h). \qquad (6)$$

Then the model concatenates $i_{1,j,t}^m$ and $h_k$ as input to the U-Net.

**Category Prompt.** In order to make the model's generated results rely more on its learned fashion matching knowledge rather than the provided prompt, the model uses a brief text description that contains only the relevant fashion categories. For example, in Figure 2, it use "A photo of a hat, with the white background" as prompt and use the CLIP text encoder [31] to encode text and obtain the Category Prompt $p_k$.

In summary, at each sampling time step $t$, the model fuses the features of the mutual condition into the noisy image and concatenates the image features from the user interaction history. These combined latent features are then fed into UNet for noise prediction.

### 4.2 Feedback Generation from Multiple Experts

From the above inference process, we obtain multiple generated fashion items $I_0 = \{i_{k,j,0}\}_{j=1}^m$. To ensure a comprehensive and objective evaluation, we designed a framework involving multiple experts to evaluate $I_0$ from three aspects: quality, compatibility, and personalization. The feedback is then utilized to fine-tune the model at each iteration. Next, we introduce three evaluation perspectives and their corresponding expert models:

**Quality.** When evaluating the quality of generated images, traditional metrics such as FID [30] and LPIPS [52] typically capture only low-level statistical features and fail to fully reflect the complex semantic information perceived by humans. Therefore, we use the multimodal large model MiniCPM [47] to score the images. We input the generated fashion items $I_0$ into the model and define the task as an image quality classification problem. Specifically, we categorize image quality into ten levels, ranging from "1-Poor Quality" to "10-Exceptional Quality," and provide the following evaluation criteria: "Consider whether the fashion elements in the image are complete and whether they conform to fashion design principles." The MiniCPM is responsible for determining which quality level the

input image $i_{k,j,0}$ belongs to and providing a response accordingly. By analyzing the image quality classification levels in MiniCPM's response, we can derive the image quality score $S_q = \{s_q^j\}_{j=1}^m$.

**Compatibility.** We follow one of the typical methods VBPR [12] to build the compatibility score prediction module. For example, given the generated fashion item $i_{1,1,0}$ and its corresponding incomplete outfit $O_1'$, we employ the pre-trained deep CNN model, i.e., ResNet-50 [11] to extract the visual features and then leverage a linear layer to transform the visual features into a shared representation space, which is formally defined as:

$$v_i = \text{ResNet}(i_{1,1,0})\mathbf{W}_1,$$
$$v_o = \text{ResNet}(O_1')\mathbf{W}_1 = \frac{1}{3}\sum_{k=2}^4 \text{ResNet}(i_{k,0})\mathbf{W}_1, \qquad (7)$$

where $v_i, v_o \in \mathbb{R}^d$ are visual representations, $\mathbf{W}_1 \in \mathbb{R}^{2048 \times d}$ is the feature transformation matrix of the linear layer, $\text{ResNet}(\cdot)$ represents the ResNet-50 network, and $d$ is the dimensionality of the visual representation. Thereafter, we can calculate the compatibility score $s_{i,o}$ via the following equation:

$$s_{i,o} = \alpha + \beta_i + \beta_o + v_i^\top v_o, \qquad (8)$$

where $\alpha, \beta_i, \beta_o$ are trainable parameters to model the bias for the global, item $i$, and incomplete outfit $O'$. We use Bayesian Personalized Ranking (BPR) [33] loss to train the model, denoted as:

$$\mathcal{L}^{\text{BPR}} = \sum_{(i,r,o)\in Q} -\log\sigma(s_{i,o} - s_{r,o}), \qquad (9)$$

where $Q = \{i,r,o\}|s_{i,0} = 1, s_{r,0} = 0\}$, $s_{i,o}$ is the ground-truth matching relation between item $i$ and incomplete outfit $O'$, $s_{i,o}$ indicates that $(i, O')$ are matched with each other. In contrast, $s_{r,o} = 0$ means the random item $r$ and $O'$ are an unmatched pair. And $\sigma(\cdot)$ represents the sigmoid function. The generated fashion items $I_0 = \{i_{k,j,0}\}_{j=1}^m$ received a set of scores $S_c = \{s_{i,o}^j\}_{j=1}^m$ after being evaluated by compatibility assessment expert.

**Personalization.** The user's interaction history with fashion items reflects their preferences. We designed a personalization evaluation expert to assess whether the generated results satisfies these preferences. Specifically, We use a pretrained CLIP image encoder to encode the generated fashion item $i_{k,j,0}$, obtaining $v_k$, and mapping it to the same latent space as the history condition $h_k$. Then, we measure the similarity between them by calculating their cosine similarity:

$$s_p = \text{CLIP\_Score}(v_k, h_k) = \frac{v_k \cdot h_k}{\| v_k \| \| h_k \|}. \qquad (10)$$

Similarly, we can obtain the personalization score $S_p = \{s_p^j\}_{j=1}^m$. Then we combine the evaluation results of multiple experts to obtain a weighted score:

$$score = \alpha_q \cdot \text{norm}(S_q) + \alpha_c \cdot \text{norm}(S_c) + \alpha_p \cdot \text{norm}(S_p), \qquad (11)$$

where the $\alpha_q, \alpha_c, \alpha_p$ is a hyperparameter that controls the weight of expert evaluation, and $\text{norm}(\cdot)$ denotes the min-max normalization. We compare the score with the threshold $t$: if $score_{i_{k,j,0}} > t$, it indicates that the generated item meets the design and personalization preferences, and it is marked as "Good" with a label of $w$. IF $score_{i_{k,j,0}} < t$, it means the generated item does not meet

the preferences, and it is marked as "Bad" with a label of $l$. Then, we construct positive-negative pairs by comparing fashion items generated within the same outfit. Using the combination formula $C(n, 2)$, we can select all possible pairs:

$$P = \left\{ (i_{k,a,0}, i_{k,b,0}) \mid 1 \le a < b \le m \right\}. \qquad (12)$$

We construct the final preference pair set $P_{pref}$ by selecting item pairs where one item matches the preference $w$ and the other does not match the preference $l$:

$$P_{\text{pref}} = \left\{ \left( i_{k,a,0}^w, i_{k,b,0}^l \right), \left( i_{k,a,0}^l, i_{k,b,0}^w \right) \right\}. \qquad (13)$$

This strategy can effectively capture the preferences within the generated fashion items, providing feedback that informs and guides the subsequent model fine-tuning process.

## 4.3 Model Fine-tuning with Direct Preference Optimization

The preference set $P_{pref}$ comprises paired fashion items, denoted as $(i^w, i^l)$, where $i^w$ represents the fashion item that aligns with multiple experts' preferences, and $i^l$ represents the item that does not. Previous work [45] has demonstrated that we can view the inference process as a multi-step Markov Decision Process (MDP). So we define $\pi(a|s)$ as the policy for taking action $a$ based on state $s$, and view the inference process from timestep $t$ to timestep $t-1$ as taking action $a_t$ from state $s_t$ to state $s_{t-1}$. In the multi-step inference process of diffusion models, we can obtain the sequence of states and actions:

$$\sigma_t = \{s_t, a_t, s_{t+1}, a_{t+1}, ..., s_T, a_T\}, 0 \le t \le T. \qquad (14)$$

Since the predicted fashion items is filled with noise in the early stages of the inference process, it is difficult for both humans and the discriminator model to judge the quality of the images. Therefore, following Reinforcement Learning (RL) [36, 37] methods, we assume that if the final inference result $i^w$ is better than $i^l$, then at any timestep during the inference process, the state $s_w$ and action $a_w$ are better than $s_l$ and $a_l$.

In Section 4.1, we saved the latent image variables from $T$ timesteps of each generated fashion item $i$ as the state $s$. Based on the positive-negative pairs $(i^w, i^l)$ obtained from the feedback, we have the states $s^w = \{i_0^w, ..., i_T^w\}$ and $s^l = \{i_0^l, ..., i_T^l\}$. As shown in Figure 2, at timestep $t \in \{1, ..., T\}$, saved latent $i_t^w$ undergoes noise prediction, guided by the conditions $(m_k, h_k, p_k)$, where the trainable and the frozen UNet respectively obtain $\epsilon_{t,\theta}^w$ and $\epsilon_{t,ref}^w$. Then we estimate the original noise-free latent variable $\hat{i}_0^w$ using the noise $\epsilon_{t,\theta}^w$ and the current latent variable $i_t^w$:

$$\hat{i}_0^w = \frac{i_t^w - \sqrt{1 - \alpha_t} \cdot \epsilon_{t,\theta}^w (i_t^w, t, c)}{\sqrt{\alpha_t}}, \qquad (15)$$

where $\alpha_t$ is a decay coefficient that we pre-calculate and store during the inference process, and $c$ represents multiple conditions. Then we can calculate the mean $\mu_{t-1}$ and variance $\sigma_t^2$ of the noise

**Table 1: Performance comparison between our method and various baselines. "Comp." and "Per." denote compatibility and personalization, respectively. Bold indicates the best results while underline denotes the second best results.**

| Dataset | iFashion | | | | | | | Polyvore-U | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Task | PFITB | | | | GOR | | | PFITB | | | | GOR | | |
| Evaluation metric | IS | IS-acc | Comp. | Per. | IS | IS-acc | Per. | IS | IS-acc | Comp. | Per. | IS | IS-acc | Per. |
| SD-v1.5 [34]* | 22.54 | 0.76 | 0.08 | 46.31 | 23.20 | 0.77 | 46.45 | 17.10 | 0.73 | 0.70 | 51.05 | 16.95 | 0.73 | 50.99 |
| SD-v2* [34] | 21.66 | 0.71 | 0.04 | 46.60 | 22.19 | 0.74 | 46.60 | 14.83 | 0.68 | 0.60 | 51.29 | 14.88 | 0.67 | 51.23 |
| SD-v1.5 [34] | 26.76 | 0.83 | 0.46 | 53.16 | 26.90 | 0.84 | 53.24 | 17.12 | 0.72 | 0.75 | 58.20 | 17.24 | 0.72 | 58.16 |
| SD-v2 [34] | 25.85 | 0.80 | 0.39 | 52.99 | 25.82 | 0.82 | 53.06 | 15.59 | 0.67 | 0.71 | 58.79 | 16.33 | 0.70 | 58.91 |
| SD-naive [34] | 25.45 | 0.80 | 0.36 | 52.95 | 25.43 | 0.81 | 52.95 | 15.45 | 0.66 | 0.73 | 59.24 | 15.48 | 0.67 | 59.12 |
| ControlNet [51] | 27.76 | 0.81 | 0.16 | 49.90 | 28.49 | 0.82 | 49.91 | 18.93 | 0.77 | 0.73 | 55.44 | _19.21_ | 0.77 | 55.40 |
| DiFashion [43] | _29.99_ | _0.90_ | _0.58_ | _55.86_ | _30.04_ | _0.90_ | _55.54_ | _19.67_ | _0.84_ | _0.80_ | _61.44_ | 18.95 | _0.83_ | _61.16_ |
| **FashionDPO(Ours)** | **33.80** | **0.91** | **0.74** | **60.39** | **32.37** | **0.91** | **59.98** | **24.14** | **0.89** | **0.83** | **64.67** | **24.93** | **0.87** | **64.79** |

distribution at time step $t - 1$:

$$\mu_{t-1} = \sqrt{\alpha_{t-1}}\hat{i}_0^w + \sqrt{1 - \alpha_{t-1}}\epsilon_{t,\theta}^w \left(i_t^w, t, c\right),$$

$$\sigma_t^2 = \frac{1 - \alpha_{t-1}}{1 - \alpha_t} \cdot \left(1 - \frac{\alpha_t}{\alpha_{t-1}}\right). \quad (16)$$

Based on the above parameters, the latent variable $i_{t-1}^w$ at timestep $t - 1$ follows a Gaussian distribution:

$$\pi_\theta \left(i_{t-1}^w \mid i_t^w, c\right) \sim \mathcal{N}\left(\mu_{t-1}, \sigma_t^2\right),$$

$$\pi_\theta \left(i_{t-1}^w \mid i_t^w, c\right) = \frac{1}{\sqrt{2\pi\sigma_t^2}} \exp\left(-\frac{\left(i_{t-1}^w - \mu_{t-1}\right)^2}{2\sigma_t^2}\right), \quad (17)$$

where $i_{t-1}^w$ refers to the latent variable that we save during the inference process at the corresponding timestep $t - 1$. The same approach can be applied to obtain other parametrized policies $\pi_\theta$ and $\pi_{ref}$. Then, based on the use of Eq.2, we derive the loss function of the fashionDPO model:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{\left(s^w, s^l\right)}\left[\log \sigma\left(\beta_w \log \frac{\pi_\theta\left(i_{t-1}^w \mid i_t^w, c\right)}{\pi_{\text{ref}}\left(i_{t-1}^w \mid i_t^w, c\right)} - \beta_l \log \frac{\pi_\theta\left(i_{t-1}^l \mid i_t^l, c\right)}{\pi_{\text{ref}}\left(i_{t-1}^l \mid i_t^l, c\right)}\right)\right], \quad (18)$$

where $\beta_w$ and $\beta_l$ denote parameters used to control preference and non-preference biases, $\sigma(\cdot)$ denotes the sigmoid function. For each preference pair $(i^w, i^l)$, fine-tuning occurs at each timestep $t \in \{T, T-1, ..., 1\}$, repeated $T - 1$ times. Subsequently, the next preference pair is drawn from the preference dataset $P_{pref}$, continuing this process until all preference data has been used for fine-tuning.

## 5 Experiments

We conduct a series of experiments on two established datasets of iFashion and Polyvore-U to evaluate our method, including both quantitative and qualitative analysis, as well as intricate model study. The experiments are designed to answer following questions:

- **RQ1:** The effectiveness of the fine-tuning framework. After multiple rounds of feedback-based fine-tuning iterations, does the model show improvements in quantitative metrics compared to DiFashion and other generative models?
- **RQ2:** During multiple iterations, how does FashionDPO perform in terms of image diversity and personalization compared to itself and other baselines?

- **RQ3:** How do data and time costs, alternative implementations of experts, and hyper-parameter adjustments affect the performance of the FashionDPO framework?

### 5.1 Experimental Settings

*5.1.1 Baselines.* For the two tasks of PFITB and GOR, we compare our model with the following baselines: **1) SD-v1.5** [34]: It's a latent space diffusion model. In the model names, with "*" indicates a pre-trained model, while without "*" indicates that the model has been fine-tuned on the fashion dataset. **2) SD-v2**: It's an upgraded version of SD-v1.5. The same naming convention. **3) SD-naive**: It's a fine-tuned model based on SD-v2, where concatenate mutual influence and history condition as condition. **4) ControlNet** [51]: It's an extension model based on SD, which controls the details of generated images by introducing additional conditional inputs. **5) DiFashion** [43]: It's the SOTA generative recommendation model based on SD-v2, which uses Classifier-Free Guidance (CFG) [13] to tightly align the control conditions of the generated fashion images.

*5.1.2 Datasets.* We follow the previous works [5, 43] and use the datasets of iFashion [1] and Polyvore-U [25], which include the required data of both fashion outfit and user-fashion item interactions. For the iFashion dataset, we select 50 common fashion categories from the dataset and filter the fashion products accordingly. Each outfit consists of four fashion products. For the Polyvore-U dataset, since it only contains basic category information such as top and bottom, we used the classifier Inception-V3 [39] fine-tuned on the iFashion dataset to perform image recognition and classification on the fashion products. We divided the dataset into training and testing sets, ensuring that each user has an interaction history with more than five outfits. The model is fine-tuned based on the feedbacks from the sampling results of training set and calculate evaluation metrics based on testing set.

*5.1.3 Evaluation Metrics.* We employ a list of quantitative evaluation metrics, covering three major evaluating perspectives: **1) Quality**: We use Inception Score (IS) to evaluate the quality of the generated images. Additionally, we use IS-Accuracy (IS-acc) to further assess whether the generated images are correct regarding its category-level semantics. **2) Compatibility (Comp.)**: We follow the previous work [43] and use the discriminator in OutfitGAN [28]
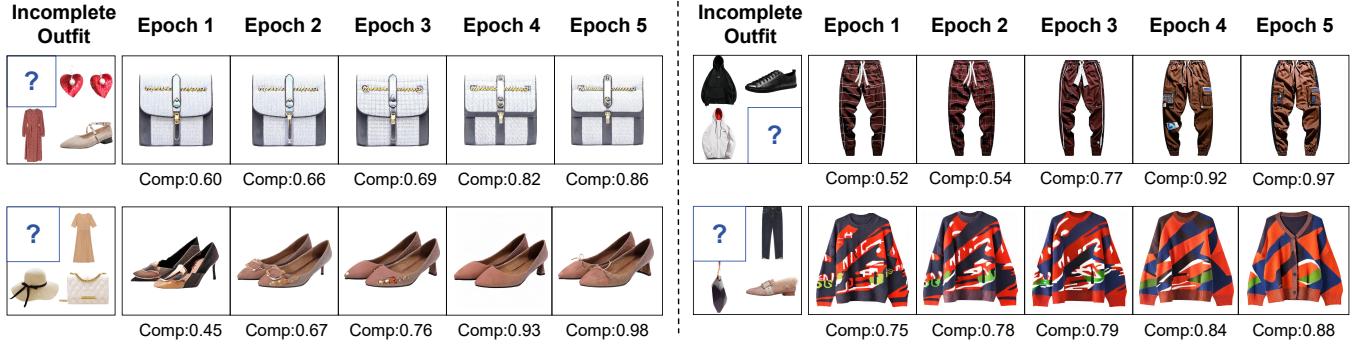
**Figure 3: Epoch-wise comparison of FashionDPO's performance across different fine-tuning epochs. As epochs increase, the compatibility metric indicates that the generated fashion items better match the incomplete outfit.**



**Figure 4: Model-wise comparison of different models' generative capabilities: PFITB task above the line, GOR task below.**

to calculate the compatibility of the generated outfit in PFITB task.
**3) Personalization (Per.)**: We use the foundation model CLIP [31] to extract the image embeddings of the items that a user has interacted in the history. Then we calculate the cosine similarity between generated fashion items and history image embeddings.

*5.1.4* **Implementation Details**. In our experiments, we load the pre-trained DiFashion model to initialize the model parameters. The training subset consists of 1,000 outfits randomly selected from the iFashion or Polyvore-U dataset. During the sampling phase, for each outfit, we generate $j = 7$ fashion items as $I_0$, with time step set to 50. At each time step, we save the latent variables and the noise predicted by diffusion model as data for fine-tuning. And in the Feedback Generation from Multiple Experts module, we set $\alpha_q = \alpha_c = \alpha_p = 1$, and threshold $t$ is the average of *score*. In the Model Fine-tuning with Direct Preference Optimization module, we perform LoRA [44] fine-tuning on the 50 saved timesteps. At each time step, we compute the loss $\mathcal{L}_{DPO}$ and update the gradients. Repeat this process until the preferences in each outfit have been learned by the model. In the next epoch, a new training subset of 1000 outfits is selected, and repeat the "sampling - feedback - fine-tuning" process. In our experiment, we fine-tuned DiFashion for five epochs to obtain the final FashionDPO model. Meanwhile, we conduct parameter search and the final model uses $\beta_w = \beta_l = 0.5$.

**Table 2: Quantitative results on different ablated models, where bold font indicates the best results while underline denotes the second best results.**

| iFashion | PFITB | | | | GOR | | |
|---|---|---|---|---|---|---|---|
| Method | IS | IS-acc | Comp. | Per. | IS | IS-acc | Per. |
| w/o Feedback | 29.46 | 0.88 | 0.59 | 55.43 | 29.30 | 0.89 | 55.89 |
| w/o Quality Expert | 30.14 | 0.88 | <u>0.70</u> | <u>59.82</u> | 29.94 | 0.90 | <u>58.41</u> |
| w/o Comp. Expert | 32.79 | <u>0.90</u> | 0.61 | 58.27 | 31.90 | 0.90 | 57.80 |
| w/o Per. Expert | <u>33.10</u> | <u>0.90</u> | 0.68 | 56.92 | <u>32.41</u> | <u>0.91</u> | 56.25 |
| **FashionDPO** | **33.80** | **0.91** | **0.74** | **60.39** | **32.37** | **0.91** | **59.98** |

## 5.2 Quantitative Results (RQ1)

In this section, we compare FashionDPO with baselines across multiple metrics and validate the effectiveness of feedback and multiple experts through ablation experiments.
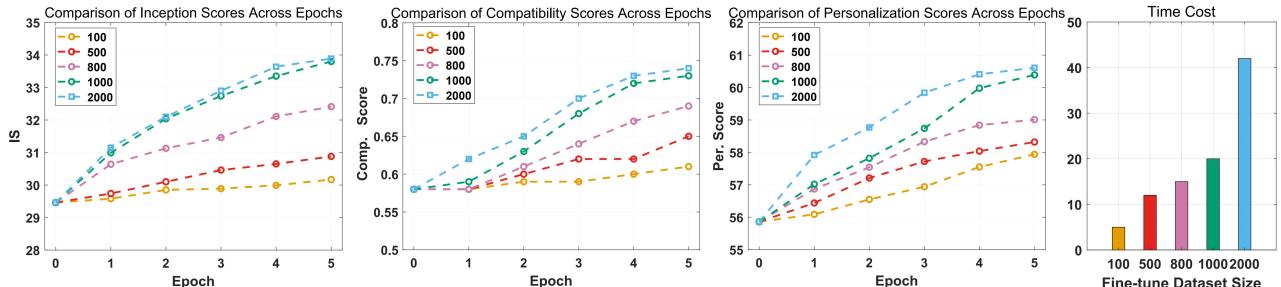
*5.2.1* **Overall Performance Comparison**. We aim to evaluate the model's performance on different datasets and tasks. The results are shown in Table 1, and we analyze them as follows:

**1) Tasks**: Our FashionDPO achieves similar scores across both tasks, indicating that the model fine-tuned on the PFITB task can smoothly transition to the GOR task. This demonstrates the generalizability of the fine-tuning framework.

**2) Models**: The results show that our FashionDPO achieves advantages over other baselines across most metrics. For the IS and IS-acc metric, the improvement of our model is significantt, indicating that our model generates more diverse results. In comparison, other baselines, such as SD-native and ControlNet, use supervised learning methods for fine-tuning, resulting in limited improvement. This indicates that our model has learned broader styling techniques and design principles. Based on the Comp. and Per. metrics, our FashionDPO demonstrates improvements over the SOTA method DiFashion. This indicates that our fine-tuning framework can effectively incorporate the knowledge of different fashion experts into the model through feedback, enabling it to learn diverse knowledge from the fashion domain.

*5.2.2* **Ablation Study**. We perform ablation experiments on the iFashion dataset, with each experiment fine-tuned for 5 epochs.
**w/o Feedback.** To verify the effectiveness of multiple experts' feedback, we removed it and fine-tuned the pre-trained DiFashion model on the iFashion dataset using LoRA.

**Figure 5: We fine-tune our model on a subset with *n* outfits, where $n \in \{100, 500, 800, 1000, 2000\}$, to explore the impact of varying datase size on model performance. Lines represent models fine-tuned on different subsets, with the x-axis as epochs and the y-axis as the evaluation metric. Bars show the time cost at per epoch (inference, feedback, fine-tuning) for different subsets.**

**Table 3: Results of expert-level human evaluation.**

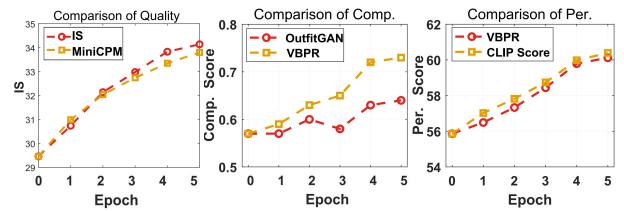| Model | D1-Style | D2-Color | D3-Fabric | D4-Variety |
|---|---|---|---|---|
| DiFashion | 2.73±0.63 | 2.74±0.64 | 2.97±0.67 | 2.91±0.47 |
| **FashionDPO(Ours)** | 4.08±0.52 | 3.87±0.37 | 3.63±0.50 | 3.22±0.43 |

**w/o Multiple Experts.** To verify the necessity of fashion experts, we conducted experiments by removing one of the three experts (*i.e.,* quality, compatibility, and personalization expert) .

The results are shown in Table 2, without feedback from multiple experts, the improvements across various metrics are minimal. By removing one of the three experts, it can be observed that without the Quality Expert has minimal impact on compatibility and personalization, while significantly affecting the IS metric. It can also be observed that removing either the Compatibility Expert or the Personalization Expert leads to significant declines in their corresponding evaluation metrics. This demonstrates that our three designed evaluation perspectives are comprehensive and objective, as the absence of any single one leads to a decline in model performance.

## 5.3 Qualitative Results (RQ2)

In addition to quantitative results, we present qualitative analysis from multiple perspectives. To verify the objectivity and effectiveness of the fine-tuning results, we invite professional fashion designers to conduct expert-level human evaluations.

*5.3.1* ***Epoch-wise Comparison****.* To validate that the model can learn fashion domain knowledge throughout the fine-tuning epochs, we present the model's performance at different epochs. As shown in Figure 3, with the increase in the number of epochs, the fashion items generated by our FashionDPO become progressively more compatible with the incomplete outfit. For example, in Epochs 1 and 2, although the generated pants and shoes roughly match the style, there are inconsistencies in materials or color tones. By Epochs 4 and 5, the generated pants and shoes align more closely with the theme of the complete outfit, and their materials and colors gradually become more consistent with the existing clothing. Notably, to facilitate the comparison of model performance across different epochs, we started from the same Gaussian noise image and used the DDIM [38] scheduler for sampling. And if we use the DDPM [17] or PNDM [22] scheduler or initialize with different noise, the results will be diverse.



**Figure 6: The impact of expert implementations. Each line corresponds to one type of expert implementation, with the x-axis as epochs and the y-axis as the evaluation metric.**

*5.3.2* ***Model-wise Comparison****.* In order to qualitative demonstrate the capability difference across various models, we present a case study by comparing the generated results of three models on two tasks. As shown in Figure 4, the results generated by our FashionDPO are more diverse. In the first row of the PFITB task results, the incomplete outfit includes a pair of black high heels and a black bag. The SD-v2 and DiFashion generate plaid shirts, which have a casual style and appear mismatched. Our FashionDPO generates light blue shirt and beige cardigan align better with the given fashion items. In the second row, the dresses generated by FashionDPO exhibit variations in style (such as floral patterns), making them more visually appealing. And in the GOR task results, FashionDPO generates a light-colored bag and brown shoes that display strong cohesion with the gray coat and skirt. While the outfits generated by SD-v2 and DiFashion are reasonable, their overall style appears somewhat inconsistent. By comparing with other baselines, it can be seen that the fashion products generated by FashionDPO adhere to fashion design principles and better meet the personalized needs of users.

*5.3.3* ***Human Evaluation by Fashion Experts****.* We invite five professional fashion designers [1] aged 18-30 to conduct expert-level double-blind human evaluations. We use a five-level scoring protocol to score 30 sets of results generated by two different models regarding to four evaluating aspects of D1-D4 (style, color, fabric, variety). The scores are discretized to five levels: very satisfied, satisfied, average, dissatisfied, and very dissatisfied, corresponding to 5 to 1 point, respectively. Finally, we collect valid scoring results from five designers for analysis.

---

[1] From Jiangnan University, School of Digital Technology & Innovation Design, 214122, Wuxi, China

**Table 4: The performance of different evaluation models on the test set. The columns of the table represent three evaluation metrics, while the rows represent the evaluation experts corresponding to different perspectives.**

| Metric | Quality | | Compatibility | | Personalization | |
|---|---|---|---|---|---|---|
| | **MiniCPM** | GPT-4 | OutfitGAN | **VBPR** | **CLIP Score** | VBPR |
| Accuracy | 89.10% | 82.50% | 70.10% | 84.80% | 87.40% | 86.60% |
| Mean | 0.56±0.15 | 0.64±0.21 | 0.47±0.14 | 0.44±0.05 | 0.53±0.04 | 0.56±0.05 |
| Pair-t | $9.29 \times 10^{-4}$ | | $6.52 \times 10^{-6}$ | | $8.31 \times 10^{-3}$ | |

As shown in Table 3, it can be seen that our FashionDPO is recognized by the majority of fashion experts, achieving higher mean scores with lower variance. Furthermore, we conduct interviews with the participating fashion designers and gathered some common subjective feedback. For instance, almost all fashion experts indicated that FashionDPO demonstrate better diversity and outfit compatibility compared to DiFashion, especially in the matching of formal wear and evening dress styles. In addition to discussions on style, the performance of individual items also catch the attention of fashion designers. For example, FashionDPO deliver satisfactory generation results in categories such as bags, shoes and dresses.
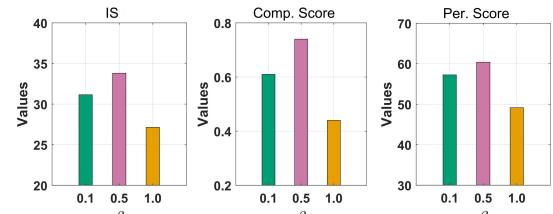
## 5.4 Model Study (RQ3)

We further study several essential properties of our model.

*5.4.1  Data and Time Cost Analysis.* We test the data and time required for fine-tuning. During iterative fine-tuning, we evaluate the model's performance differences across different epochs using three major evaluation perspectives: Quality - IS, Compatibility - Comp. Score, and Personalization - Per. Score. As shown in Figure 5, we present the results fine-tuned on five subsets of different sizes. As the size of the subset increases, the model can learn fashion knowledge within less epochs, while the time required for fine-tuning still increases. Notably, when the number of epochs is set to five, the performance of the subset with 1,000 is similar to that with 2,000. However, the time doubles for the subset with 2,000 outfits. Therefore, we choose 1,000 outfits as the default fine-tuning subset.

*5.4.2  Alternative Implementations of Experts.* To validate whether our method is sufficiently generalizable for expert implementations, we replace the expert model in each of the three evaluation perspectives and fine-tune for five epochs on the same data subset. The results are shown in Figure 6, where we present the performance for different expert implementations, with the x-axis as epochs and the y-axis as the evaluation metric. We can see that replacing the quality and personalization experts has little impact, whereas the compatibility expert has a greater influence. Further investigation revealed that OutfitGAN's scores are relatively extreme, indicating poor generalization performance. This demonstrates that improving the quality of the expert can further enhance the model's performance.

We further tested the quality of the multi-expert feedback module to ensure that the scores given by these experts are comprehensive and objective. As shown in Table 4, in terms of Quality, MiniCPM outperforms GPT-4 in accuracy. This is because GPT-4 has a weaker ability to distinguish between positive and negative fashion items



**Figure 7: Effects of the hyper-parameter $\beta$ in controlling the deviation. Bars show the scores on the evaluation metric after fine-tuning for five epochs using different hyper-parameters.**

compared to MiniCPM, as evidenced by GPT-4 assigning relatively high scores to negative items as well. In terms of Compatibility, the Accuracy of VBPR is significantly higher than that of Outfit-GAN. This is because Outfit-GAN's scoring tends to be more extreme, resulting in a lack of distinction. In terms of Personalization, both CLIP-Score and VBPR perform similarly in accuracy, and their efficiency in model fine-tuning is also comparable. From the Pair-t test, it is evident that the p-values between the models for all three evaluation aspects are below the set significance level (0.05), indicating that there are differences between the models. It further demonstrates the versatility of the FashionDPO framework, showing that the experts within it are interchangeable, and that stronger experts lead to better fine-tuning performance.

*5.4.3  Hyper-parameter Analysis.* The hyper-parameter $\beta = \beta_w = \beta_l$ is used to control the preference and non-preference biases. When $\beta$ increases, it amplifies the differences between preferences and accelerates the model's convergence. We fine-tune FashionDPO for five epochs with different $\beta$ and test it on three major evaluation perspectives. The performance regarding various $\beta$ is shown in Figure 7. We can see that when $\beta$ increases from 0.1 to 1.0, the performance experience a pattern of first growing then dropping. When $\beta$ is too small, it becomes difficult to distinguish the differences between preferred and non-preferred data, thereby affecting the model's efficiency in learning preferences. When is $\beta$ too large, it overly amplifies the preference differences, which can cause the model to fall into a local optimum. Therefore, we chose hyper-parameter $\beta = 0.5$.

## 6  Conclusion And Future Work

In summary, we identified the limitations of lacking diversity and supervised learning paradigm in GOR models, and we propose to leverage AI feedback to address the problem. To address the challenge of the design of AI evaluation models and the mechanisms for updating feedback, we proposed a novel framework FashionDPO, which fine-tunes the fashion outfit generation model using multiple automated experts' feedbacks.

For future work, we will introduce more expert perspectives and explore the interactions among different experts. Second, regarding the feedback mechanism, we will incorporate the intensity of preferences to refine the distinctions between preferences and non-preferences. This will enhance our framework's ability to learn expert knowledge. Furthermore, we will explore additional feedback mechanisms, such as integrating verbal feedback from LLMs into the generation model.

# References

[1] Wen Chen, Pipei Huang, Jiaming Xu, Xin Guo, Cheng Guo, Fei Sun, Chao Li, Andreas Pfadler, Huan Zhao, and Binqiang Zhao. 2019. POG: Personalized Outfit Generation for Fashion Recommendation at Alibaba iFashion. In *KDD*. ACM, 2662–2670.

[2] Yujuan Ding, Zhihui Lai, P. Y. Mok, and Tat-Seng Chua. 2024. Computational Technologies for Fashion Recommendation: A Survey. *ACM Comput. Surv.* 56, 5 (2024), 121:1–121:45.

[3] Yujuan Ding, Yunshan Ma, Wenqi Fan, Yige Yao, Tat-Seng Chua, and Qing Li. 2024. FashionReGen: LLM-Empowered Fashion Report Generation. In *WWW (Companion Volume)*. ACM, 991–994.

[4] Yujuan Ding, Yunshan Ma, Lizi Liao, Wai Keung Wong, and Tat-Seng Chua. 2022. Leveraging Multiple Relations for Fashion Trend Forecasting Based on Social Media. *IEEE Trans. Multim.* 24 (2022), 2287–2299.

[5] Yujuan Ding, P. Y. Mok, Yunshan Ma, and Yi Bin. 2023. Personalized fashion outfit generation with user coordination preference learning. *Inf. Process. Manag.* 60, 5 (2023), 103434.

[6] Xue Dong, Xuemeng Song, Fuli Feng, Peiguang Jing, Xin-Shun Xu, and Liqiang Nie. 2019. Personalized Capsule Wardrobe Creation with Garment and User Modeling. In *ACM Multimedia*. ACM, 302–310.

[7] Xiaoyu Du, Kun Qian, Yunshan Ma, and Xinguang Xiang. 2023. Enhancing item-level bundle representation for bundle recommendation. *ACM Transactions on Recommender Systems* (2023).

[8] Junhong Gou, Siyu Sun, Jianfu Zhang, Jianlou Si, Chen Qian, and Liqing Zhang. 2023. Taming the Power of Diffusion Models for High-Quality Virtual Try-On with Appearance Flow. In *ACM Multimedia*. ACM, 7599–7607.

[9] Xintong Han, Zuxuan Wu, Yu-Gang Jiang, and Larry S. Davis. 2017. Learning Fashion Compatibility with Bidirectional LSTMs. In *ACM Multimedia*. ACM, 1078–1086.

[10] Xintong Han, Zuxuan Wu, Zhe Wu, Ruichi Yu, and Larry S. Davis. 2018. VITON: An Image-Based Virtual Try-On Network. In *CVPR*. Computer Vision Foundation / IEEE Computer Society, 7543–7552.

[11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep Residual Learning for Image Recognition. *CoRR* abs/1512.03385 (2015).

[12] Ruining He and Julian J. McAuley. 2016. VBPR: Visual Bayesian Personalized Ranking from Implicit Feedback. In *AAAI*. AAAI Press, 144–150.

[13] Jonathan Ho and Tim Salimans. 2022. Classifier-Free Diffusion Guidance. *CoRR* abs/2207.12598 (2022).

[14] Qihan Huang, Long Chan, Jinlong Liu, Wanggui He, Hao Jiang, Mingli Song, and Jie Song. 2024. PatchDPO: Patch-level DPO for Finetuning-free Personalized Image Generation. *CoRR* abs/2412.03177 (2024).

[15] Cong Phuoc Huynh, Arri Ciptadi, Ambrish Tyagi, and Amit Agrawal. 2018. CRAFT: Complementary Recommendations Using Adversarial Feature Transformer. *CoRR* abs/1804.10871 (2018).

[16] Jeongho Kim, Gyojung Gu, Minho Park, Sunghyun Park, and Jaegul Choo. 2023. StableVITON: Learning Semantic Correspondence with Latent Diffusion Model for Virtual Try-On. *CoRR* abs/2312.01725 (2023).

[17] Hyeon-Ju Lee and Seok-Jun Buu. 2024. Deep Generative Replay With Denoising Diffusion Probabilistic Models for Continual Learning in Audio Classification. *IEEE Access* 12 (2024), 134714–134727.

[18] Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. 2023. RLAIF: Scaling Reinforcement Learning from Human Feedback with AI Feedback. *CoRR* abs/2309.00267 (2023).

[19] Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. 2024. RLAIF vs. RLHF: Scaling Reinforcement Learning from Human Feedback with AI Feedback. In *ICML*. OpenReview.net.

[20] Xingchen Li, Xiang Wang, Xiangnan He, Long Chen, Jun Xiao, and Tat-Seng Chua. 2020. Hierarchical Fashion Graph Network for Personalized Outfit Recommendation. In *SIGIR*. ACM, 159–168.

[21] Zhanhao Liang, Yuhui Yuan, Shuyang Gu, Bohan Chen, Tiankai Hang, Ji Li, and Liang Zheng. 2024. Step-aware Preference Optimization: Aligning Preference with Denoising Performance at Each Step. *CoRR* abs/2406.04314 (2024).

[22] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. 2022. Pseudo Numerical Methods for Diffusion Models on Manifolds. In *ICLR*. OpenReview.net.

[23] Xiaohao Liu, Jie Wu, Zhulin Tao, Yunshan Ma, Yinwei Wei, and Tat-Seng Chua. 2025. Fine-tuning Multimodal Large Language Models for Product Bundling. In *KDD (1)*. ACM, 848–858.

[24] Zhi Lu, Yang Hu, Yan Chen, and Bing Zeng. 2021. Personalized Outfit Recommendation With Learnable Anchors. In *CVPR*. Computer Vision Foundation / IEEE, 12722–12731.

[25] Zhi Lu, Yang Hu, Yunchao Jiang, Yan Chen, and Bing Zeng. 2019. Learning Binary Code for Personalized Fashion Recommendation. In *CVPR*. Computer Vision Foundation / IEEE, 10562–10570.

[26] Yunshan Ma, Yujuan Ding, Xun Yang, Lizi Liao, Wai Keung Wong, and Tat-Seng Chua. 2020. Knowledge Enhanced Neural Fashion Trend Forecasting. In *ICMR*.

[27] Yunshan Ma, Yingzhi He, Xiang Wang, Yinwei Wei, Xiaoyu Du, Yuyangzi Fu, and Tat-Seng Chua. 2024. MultiCBR: Multi-view Contrastive Learning for Bundle Recommendation. *ACM Trans. Inf. Syst.* 42, 4 (2024), 100:1–100:23.

[28] Maryam Moosaei, Yusan Lin, Ablaikhan Akhazhanov, Huiyuan Chen, Fei Wang, and Hao Yang. 2022. OutfitGAN: Learning Compatible Items for Generative Fashion Outfits. In *CVPR Workshops*. IEEE, 2272–2276.

[29] Sanghyeon Na, Yonggyu Kim, and Hyunjoon Lee. 2024. Boost Your Own Human Image Generation Model via Direct Preference Optimization with AI Feedback. *CoRR* abs/2405.20216 (2024).

[30] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. 2022. On Aliased Resizing and Surprising Subtleties in GAN Evaluation. In *CVPR*. IEEE, 11400–11410.

[31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *ICML (Proceedings of Machine Learning Research, Vol. 139)*. PMLR, 8748–8763.

[32] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. In *NeurIPS*.

[33] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian Personalized Ranking from Implicit Feedback. In *UAI*. AUAI Press, 452–461.

[34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis with Latent Diffusion Models. In *CVPR*. IEEE, 10674–10685.

[35] Yong-Siang Shih, Kai-Yueh Chang, Hsuan-Tien Lin, and Min Sun. 2018. Compatibility Family Learning for Item Recommendation and Generation. In *AAAI*. AAAI Press, 2403–2410.

[36] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Vedavyas Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy P. Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. 2016. Mastering the game of Go with deep neural networks and tree search. *Nat.* 529, 7587 (2016), 484–489.

[37] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy P. Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. 2017. Mastering the game of Go without human knowledge. *Nat.* 550, 7676 (2017), 354–359.

[38] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2021. Denoising Diffusion Implicit Models. In *ICLR*. OpenReview.net.

[39] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2016. Rethinking the Inception Architecture for Computer Vision. In *CVPR*. IEEE Computer Society, 2818–2826.

[40] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. 2023. Diffusion Model Alignment Using Direct Preference Optimization. *CoRR* abs/2311.12908 (2023).

[41] Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. 2023. ELITE: Encoding Visual Concepts into Textual Embeddings for Customized Text-to-Image Generation. In *ICCV*. IEEE, 15897–15907.

[42] Zhenyu Xie, Zaiyu Huang, Xin Dong, Fuwei Zhao, Haoye Dong, Xijin Zhang, Feida Zhu, and Xiaodan Liang. 2023. GP-VTON: Towards General Purpose Virtual Try-On via Collaborative Local-Flow Global-Parsing Learning. In *CVPR*. IEEE, 23550–23559.

[43] Yiyan Xu, Wenjie Wang, Fuli Feng, Yunshan Ma, Jizhi Zhang, and Xiangnan He. 2024. Diffusion Models for Generative Outfit Recommendation. In *SIGIR*. ACM, 1350–1359.

[44] Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhengsu Chen, Xiaopeng Zhang, and Qi Tian. 2024. QA-LoRA: Quantization-Aware Low-Rank Adaptation of Large Language Models. In *ICLR*. OpenReview.net.

[45] Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiaxin Chen, Qimai Li, Weihan Shen, Xiaolong Zhou, and Xiu Li. 2023. Using Human Feedback to Fine-tune Diffusion Models without Any Reward Model. *CoRR* abs/2311.13231 (2023).

[46] Zilin Yang, Zhuo Su, Yang Yang, and Ge Lin. 2018. From recommendation to generation: A novel fashion clothing advising framework. In *2018 7th International Conference on Digital Home (ICDH)*. IEEE, 180–186.

[47] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. 2024. MiniCPM-V: A GPT-4V Level MLLM on Your Phone. *arXiv preprint arXiv:2408.01800* (2024).

[48] Mingzhe Yu, Yunshan Ma, Lei Wu, Kai Cheng, Xue Li, Lei Meng, and Tat-Seng Chua. 2024. Smart Fitting Room: A One-stop Framework for Matching-aware Virtual Try-On. In *ICMR*. ACM, 184–192.

[49] Yu Zeng, Vishal M. Patel, Haochen Wang, Xun Huang, Ting-Chun Wang, Ming-Yu Liu, and Yogesh Balaji. 2024. JeDi: Joint-Image Diffusion Models for Finetuning-Free Personalized Text-to-Image Generation. In *CVPR*. IEEE, 6786–6795.

[50] Huijing Zhan, Jie Lin, Kenan Emir Ak, Boxin Shi, Ling-Yu Duan, and Alex C. Kot. 2022. $Aˆ3$-FKG: Attentive Attribute-Aware Fashion Knowledge Graph for Outfit Preference Prediction. *IEEE Trans. Multim.* 24 (2022), 819–831.

[51] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023. Adding Conditional Control to Text-to-Image Diffusion Models. In *ICCV*. IEEE, 3813–3824.

[52] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*. Computer Vision Foundation / IEEE Computer Society, 586–595.