大学共同利用機関法人 情報・システム研究機構
国立情報学研究所
National Institute of Informatics
NII

東京大学
THE UNIVERSITY OF TOKYO

# Pose aware Outfit Transfer between Unpaired in-the-wild images

Donnaphat Trakulwaranont, Marc A. Kastner, Shin'ichi Satoh

✉ eiam@nii.ac.jp

## Introduction

- Fashion is one of highest revenue industries in the world.
- Online fashion shopping has some obstacles: Hard to judge a product's look on oneself

**Motivation**

- However, outfit transfer method has various existing problems:
  - Usually trained on studio photography (No backgrounds)
  - Need for paired data
- We want to handle noisy backgrounds such as in street photography



Studio photography + Paired data

Street photography (Noisy background)

## Proposed method

Our method has 2 main modules

1. Warp mockup target clothing image generation
   - It is used for generate target clothing image close to target model body shape.

Feature extractor:
$$\begin{cases} f_{P_i} = F_{e_A}(P_i) \\ f_{P_j} = F_{e_B}(P_j, BP_j) \end{cases}$$

Correlation: $c(i,j) = f_{P_i}{}^T f_{P_j}$



Transformation parameter: $\theta(i,j) = Regression(c(i,j))$ , Transform image: $W_{P_i}^j = \mathcal{T}_\theta(P_i)$

2. Pose-aware Outfit transfer
   - It is used to generate image of target person wear specific clothing, and the reconstruction part is introduced to train this module as we have real pair data.

Outfit Transfer Image: $O_j^f = PATN(B_j, W_{P_i}^j, BP_j)$ Reconstructed original Image: $O_i^f = PATN(B_i, W_{P_j}^i, BP_i)$

## Training

- For warp mockup target clothing image generation
  - use L1-Loss to train

$$\mathcal{L}(\theta, \theta_{GT}) = \frac{1}{N}\sum_{i=1}^{N}\|G_i' - G_i''\|$$

- For pose-aware outfit transfer
  - use GAN-Loss, combination L1-loss and reconstruction loss to train

$$\mathcal{L} = \mathcal{L}_{GAN} + \mathcal{L}_{combL1} + \mathcal{L}_{recon}$$

$$\mathcal{L}_{recon} = \left\|\Psi_k(O_i) - \Psi_k(O_i^f)\right\|_1$$

$$\mathcal{L}_{GAN} = \mathbb{E}_{BP_j \in \mathcal{P}_r, (O_i, O_j) \in X}\left\{\log\left[D_A(O_i, O_j) \cdot D_S(BP_j, O_j)\right]\right\}$$
$$+ \mathbb{E}_{BP_j \in \mathcal{P}_r, O_j \in X, O_j^f \in \tilde{X}}\left\{\log\left[\left(1 - D_A(O_j, O_j^f)\right) \cdot \left(1 - D_S(BP_j, O_j^f)\right)\right]\right\}$$

$$\mathcal{L}_1 = \left\|W_{P_i}^j - (O_j^f \otimes PM_j)\right\|_1 + \left\|B_j - (O_j^f \otimes BM_j)\right\|_1$$

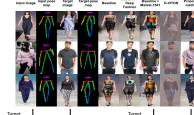$$\mathcal{L}_{perL1} = \left\|\Psi_k(W_{P_i}^j) - \Psi_k(O_j^f \otimes PM_j)\right\|_1$$
$$+ \left\|\Psi_k(B_j) - \Psi_k(O_j^f \otimes BM_j)\right\|_1$$

$$\mathcal{L}_{combL1} = \lambda_1 \mathcal{L}_1 + \lambda_2 \mathcal{L}_{perL1}$$

## Evaluation

Quantitative results

- Using Structural similarity (SSIM) and Inception score (IS)
- Promising results:
  For SSIM and mask-SSIM,
  around 54% and 45% improvement,
  while for IS around 25% improvement

| Model | SSIM | IS | Mask-SSIM | Mask-IS |
|---|---|---|---|---|
| Baseline | 0.302 | 4.073 | 0.591 | 4.444 |
| Baseline + DeepFashion | 0.282 | 4.073 | 0.580 | 4.562 |
| Baseline + Market-1501 | 0.256 | 3.912 | 0.565 | 4.064 |
| O-VTON | 0.253 | 2.648 | 0.320 | 3.292 |
| **Proposed method** | 0.467 | 5.096 | 0.860 | 3.994 |

Qualitative results

- Confirms a perception much closer to the expectation
- We can accurately preserve background information while still being able to correctly transfer the outfit

Real scenerio

- Test on myself
  - Outfit images from fashion shopping website
  - My image from webcam or real street photograph