# Learning to Wear: Details-Preserved Virtual Try-on via Disentangling Clothes and Wearer

Sangho Lee[1], Seoyoung Lee[1], Joonseok Lee[1,2]
[1]Seoul National University  [2]Google Research

## Introduction

### Virtual Try-on

- Synthesizing a realistic image of a person wearing the given clothing



Input (Model + Target clothes)    Output

- Need to disentangle human and clothes & different types of clothes data
- Should be generalizable to various human poses and body shapes

## Motivation

### Inability to reflect details of the target clothing

- Detailed characteristics of of the target clothing (*e.g.*, shape of neckline and sleeves) are not retained.
- Output often reveals characteristics of clothes in the reference image.

### Limited understanding in 3D semantics of wearing clothes

- Parts that should not be seen when worn (*e.g.*, inner side of shirt neckline) are still visible when clothing is worn.
- Overall struggle in synthesizing well-fitted images indicate a weak generalizability of the models.
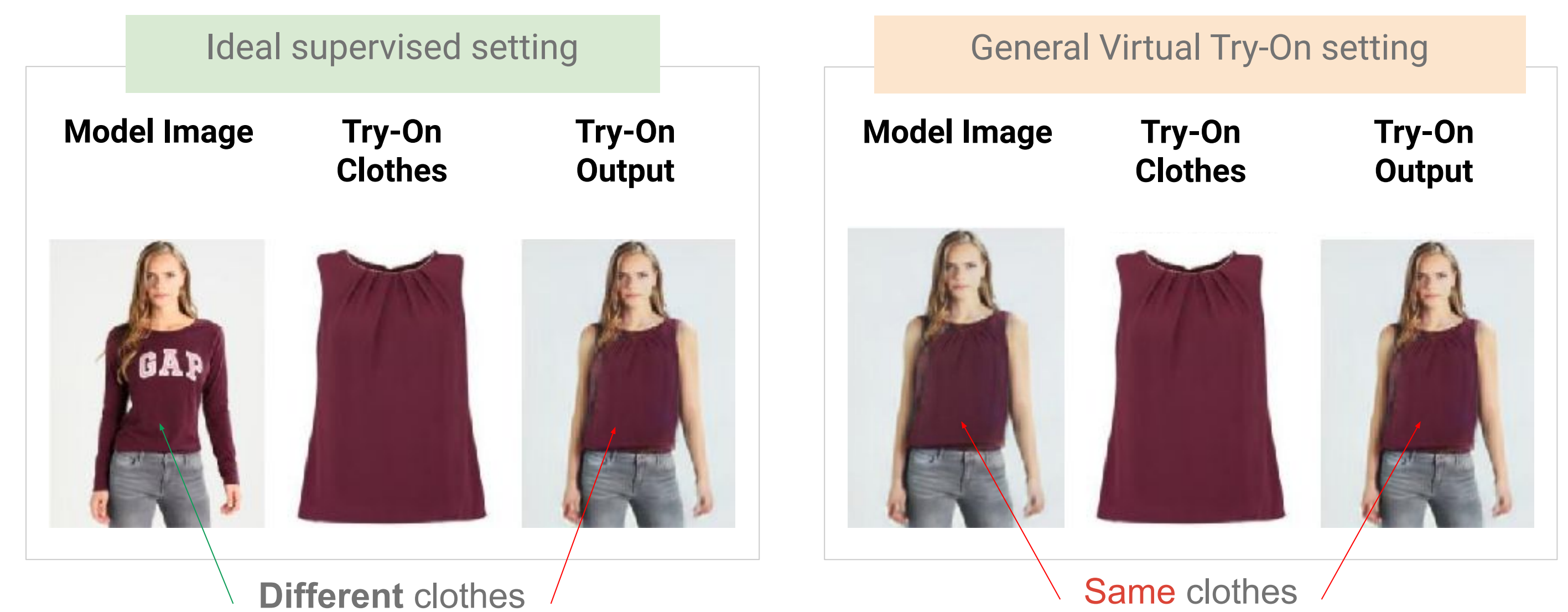
→ *Disentangle clothes & wearer and target & source clothes*



ACGPN (2020 CVPR)    PFAFN (2021 CVPR)    VITON-HD (2021 CVPR)

Reference  Target clothes  Try-on result

Reference  Target clothes  CP-VTON+  ACGPN  PFAFN

## Evaluation

### Flaw of existing metrics based on entire images



Baseline    Ours

SSIM(↑)
**0.800**    0.716

CP-VTON+  ACGPN  PFAFN  **Ours**

Ground Truth

Variance(↓)  333.8  349.0  621.6  285.2

- Misalignment between qualitative and quantitative evaluation
- Considers even parts irrelevant for virtual try-on (*e.g.*, background).

### Novel approach to evaluation metrics

- Focus on areas relevant to the task (*i.e.*, major body keypoints) for better alignment with human perception.

$$\mathtt{Metric}_\epsilon^p(I) = \frac{1}{k}\sum_{i=1}^{k}\mathtt{Metric}^{\mathrm{all}}\left(I\left[x_i - \frac{\epsilon}{2} : x_i + \frac{\epsilon}{2}, y_i - \frac{\epsilon}{2} : y_i + \frac{\epsilon}{2}\right]\right)$$
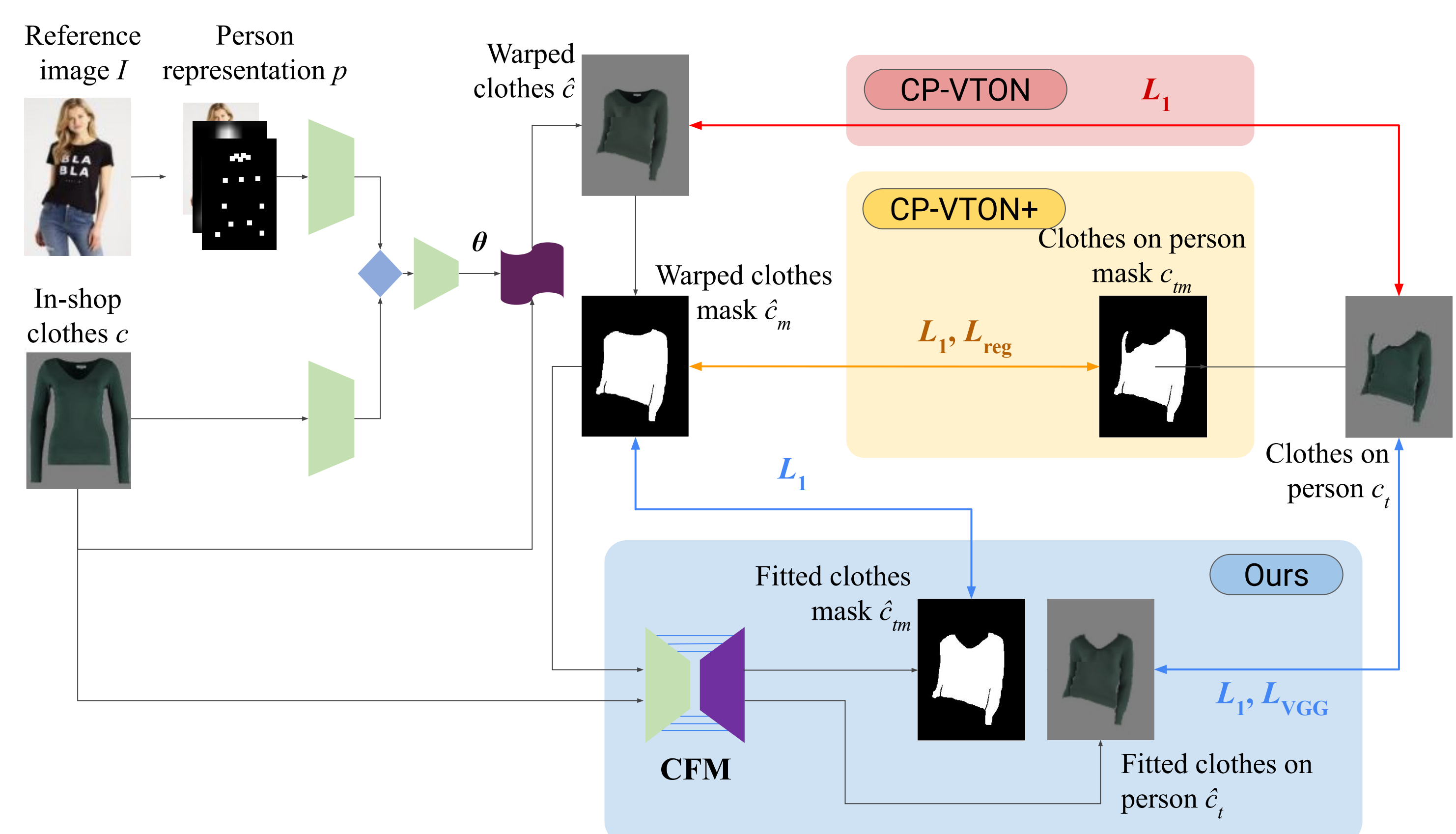


## Core Issue

### Lack of paired datasets

- Impossible to have two photos of a model with just different clothes on.
- Usually, models are trained to wear the same clothes that is already worn.



Ideal supervised setting

Model Image    Try-On Clothes    Try-On Output

**Different** clothes

General Virtual Try-On setting

Model Image    Try-On Clothes    Try-On Output

Same clothes

## Method

### Clothes Fitting Module (CFM): Learning to Wear

- Inserted between Geometric Warping Module (GWM) and Try-on Module.
- Allows GWM to perceive clothes as the source clothes in the reference.
- Allows CFM to perceive clothes as the target ground truth.
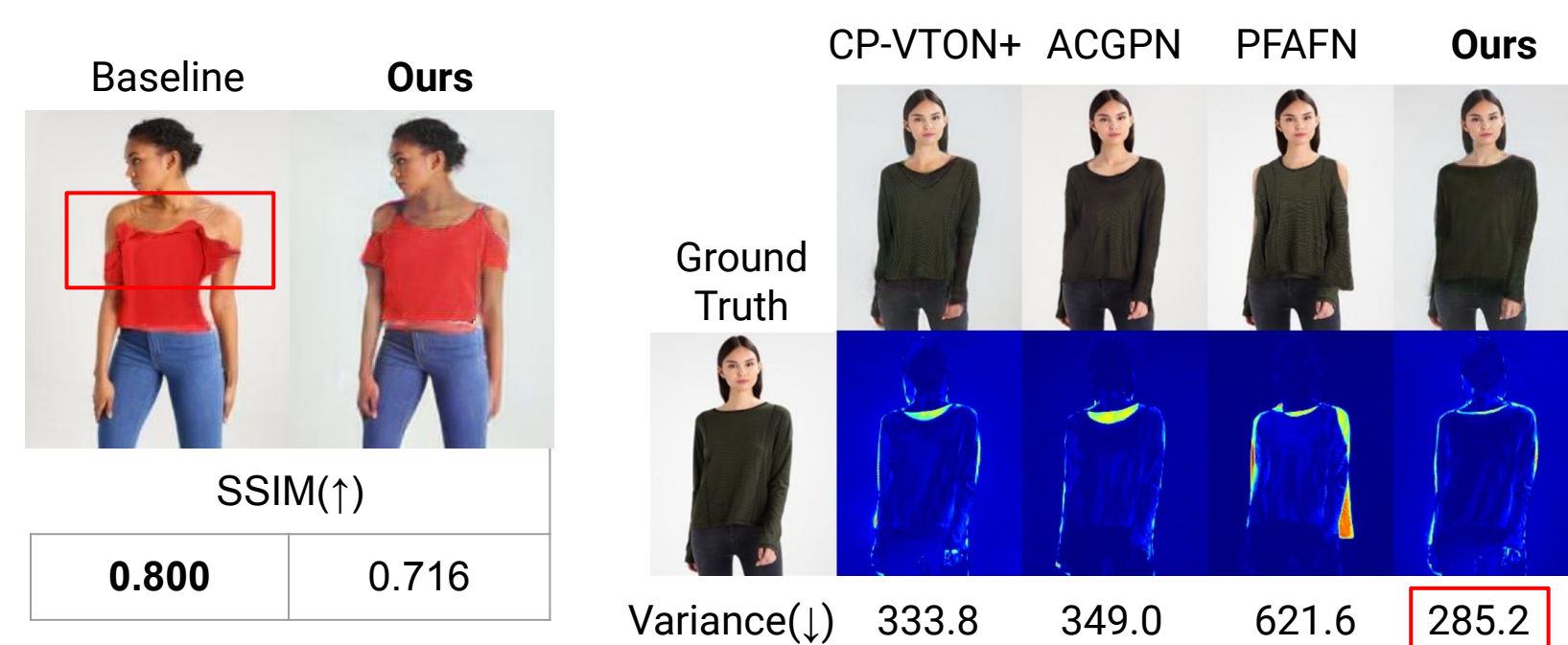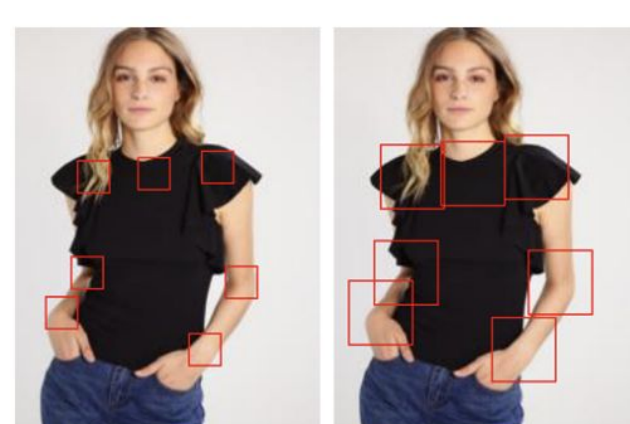- *Successfully disentangles source & target clothes and clothes & wearer*



Reference image $I$    Person representation $p$    Warped clothes $\hat{c}$

In-shop clothes $c$

$\theta$

Warped clothes mask $\hat{c}_m$

CP-VTON    $L_1$

CP-VTON+    Clothes on person mask $c_{tm}$

$L_1, L_{\mathrm{reg}}$

Clothes on person $c_t$

$L_1$

Fitted clothes mask $\hat{c}_{tm}$    Ours

CFM

$L_1, L_{\mathrm{VGG}}$

Fitted clothes on person $\hat{c}_t$

## Results

- Retain properties of target clothes, disentangled from the reference image
- Generalizable to various designs, as well as body shapes and poses.



Reference person  Target clothes  CP-VTON+  ACGPN  PFAFN  **Ours**    Reference person  Target clothes  CP-VTON+  ACGPN  PFAFN  **Ours**

## Take-home Messages

- Previous virtual try-on models learned entangled representations that lack generalizability due to the lack of paired datasets.
- With CFM, we disentangle important factors of virtual try-on and detour the inherent limitation in data.
- Our patch-based evaluation metrics better correspond to qualitative results.