

# Weakly Supervised High-Fidelity Clothing Model Generation

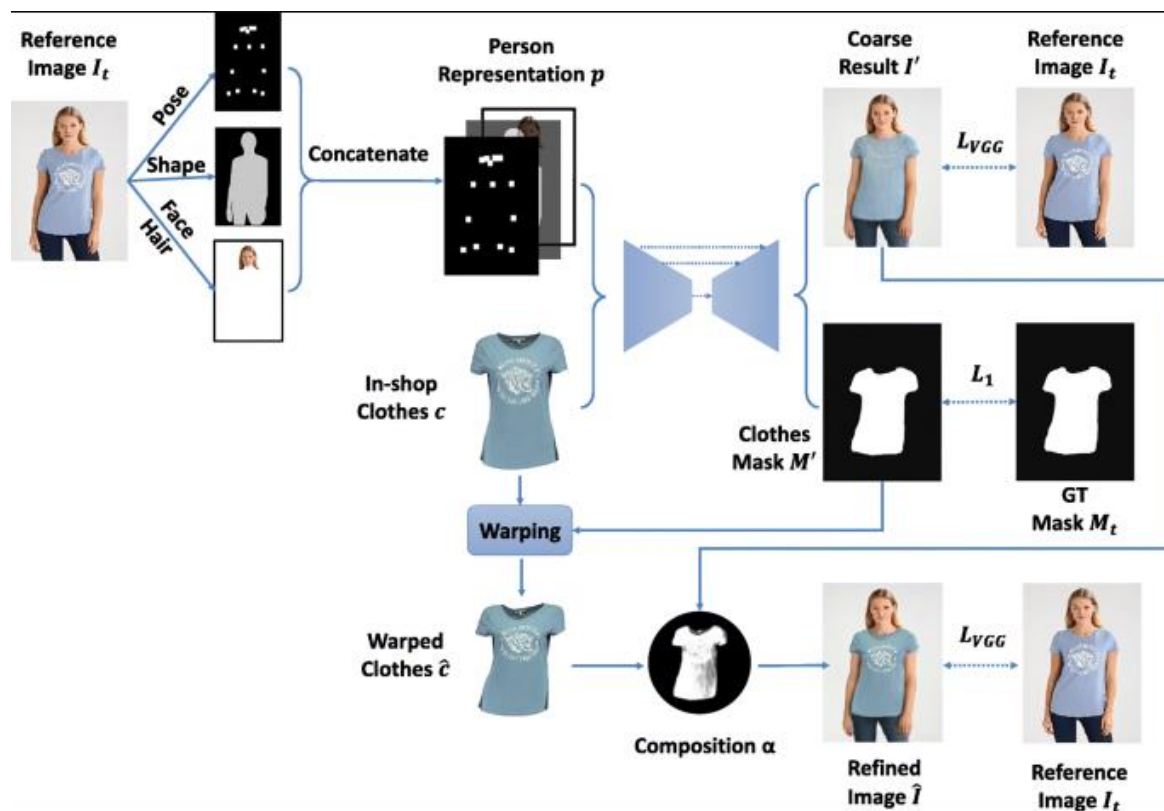
Ruili Feng, Cheng Ma, Chengji Shen, Xin Gao, Zhenjiang Liu, Xiaobo Li, Kairi Ou, Zhengjun Zha

University of Science and Technology of China   Zhejiang University   Alibaba Group

**CVPR 2022**

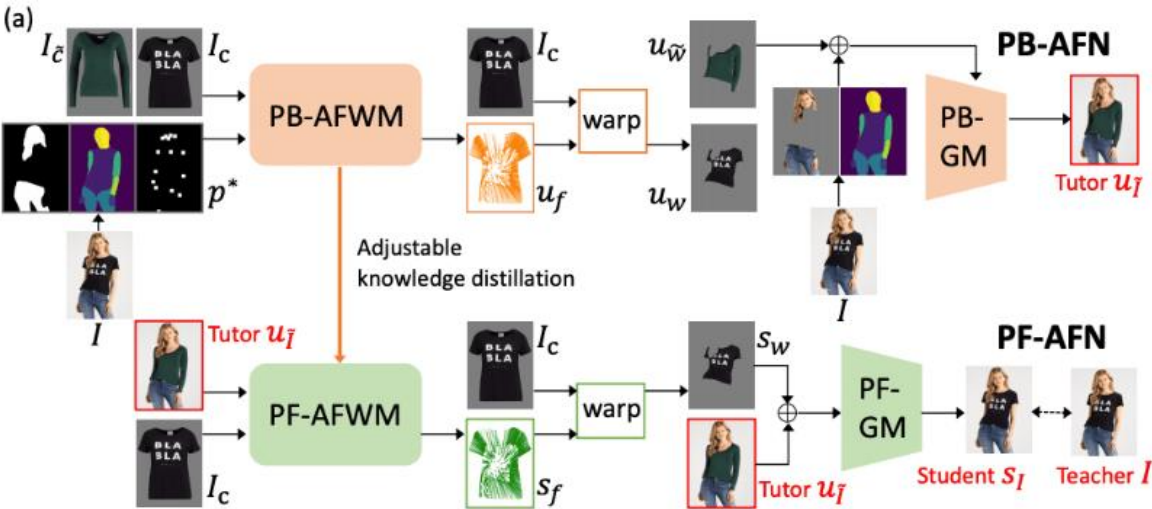


# Virtual Try-on



- VITON needs paired data
- Supervised learning

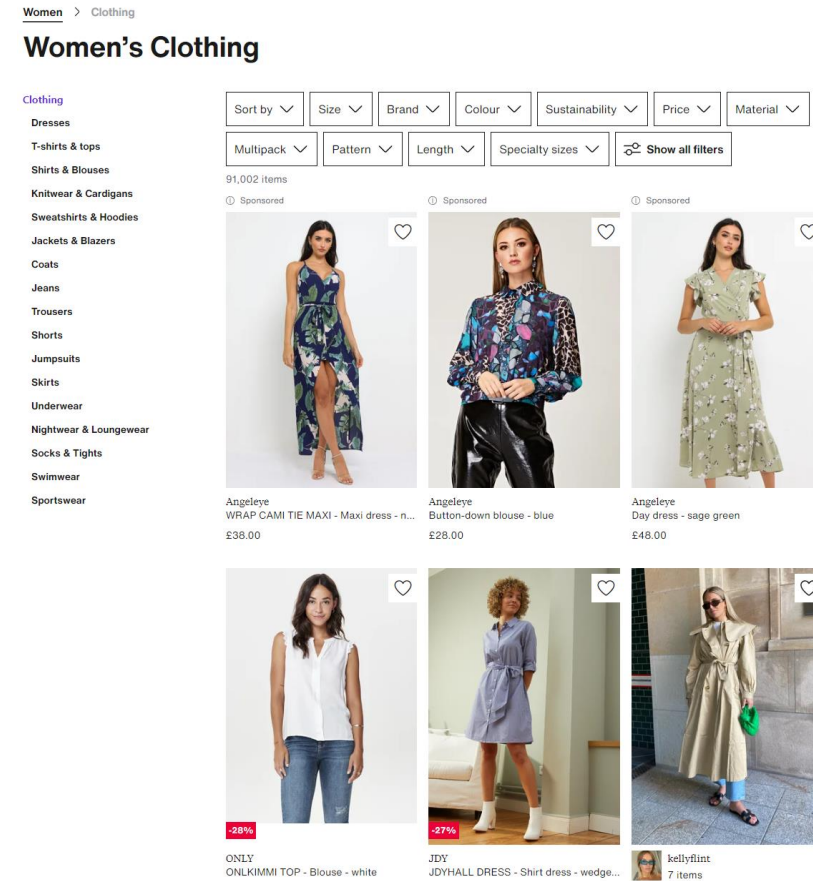
# Virtual Try-on



- “Parser-based” or “Parser-free”
- Warping + Synthesis

# Motivation

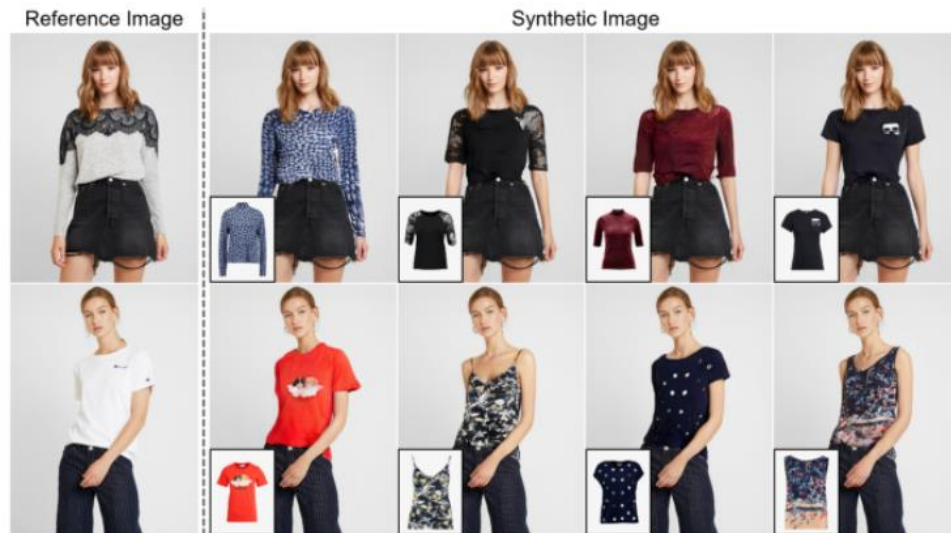
- Proprietary model images are expensive!
- Training VITON model with large dataset -> high-fidelity
- Weakly supervised learning



Zalando website

# Task setting

- Conventional VITON setting



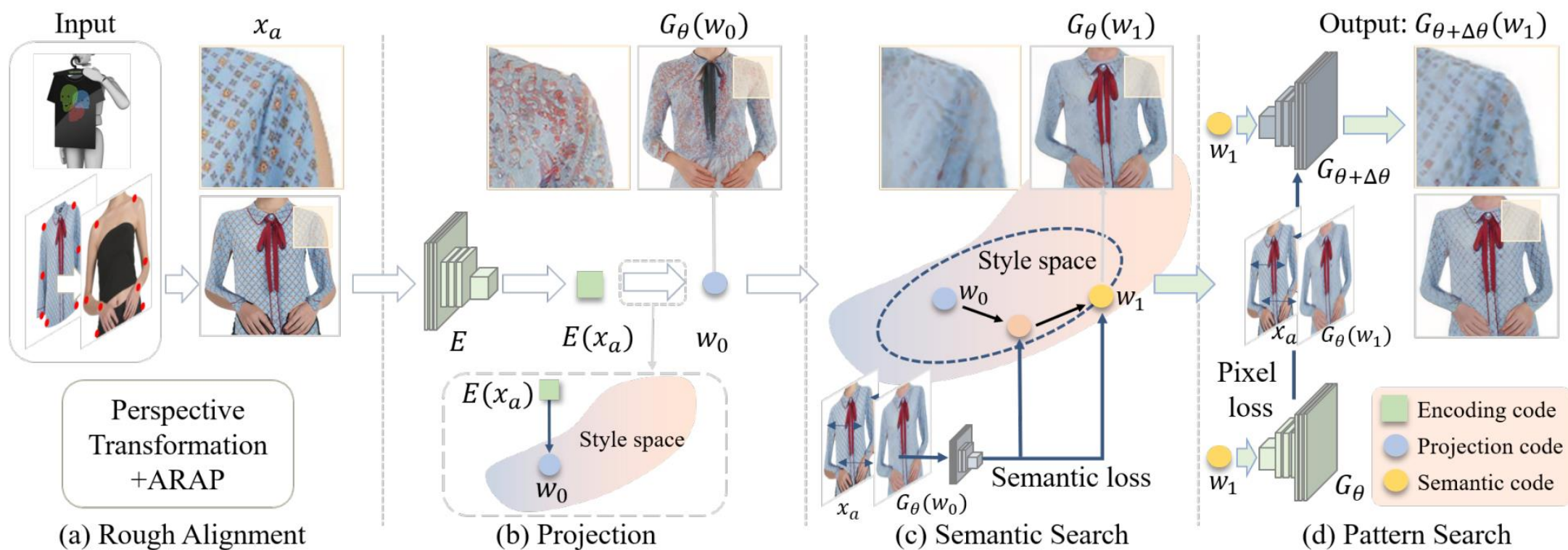
- Proposed setting



- The procedure of people predicting how they will look like while picking clothes
- Commercial Model Image dataset (CMI) – 2,348 images of models on underwear or sleeveless



# Overview



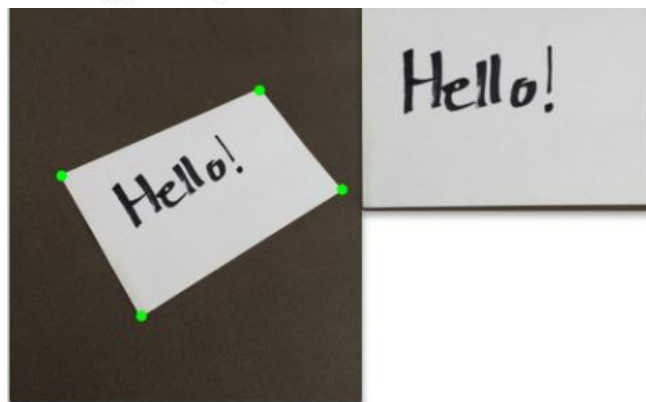
# Rough Alignment



(a) Perspective transformation

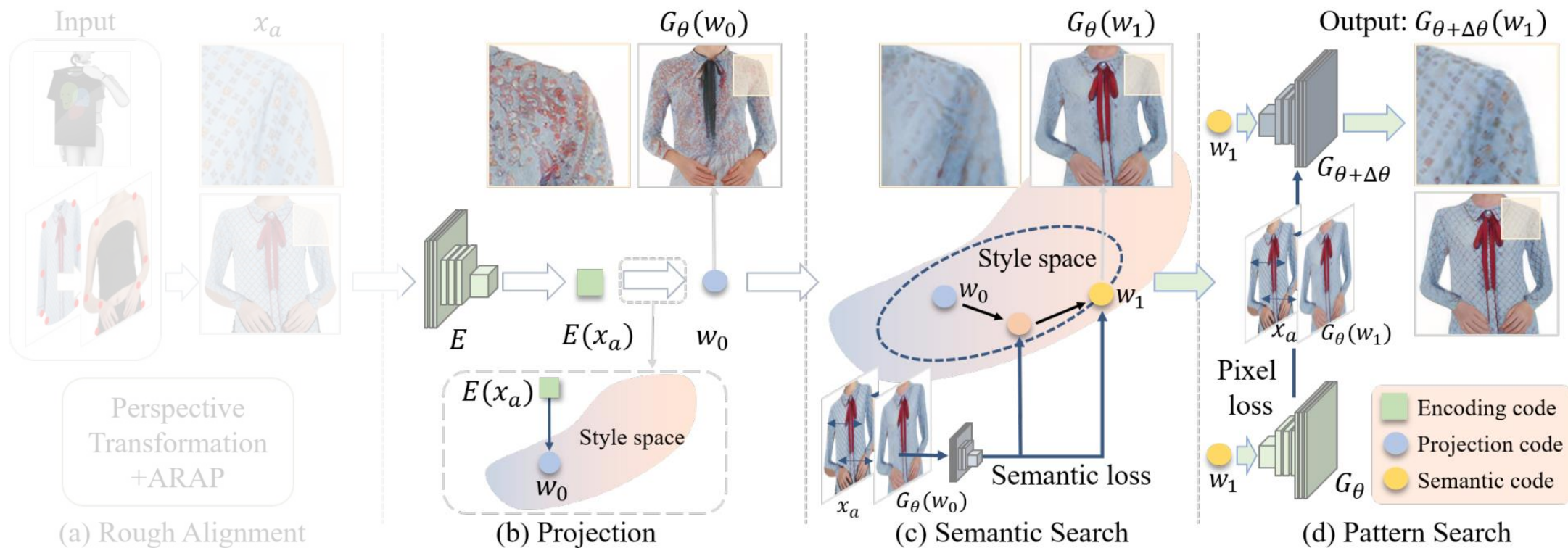


(b) ARAP



- Get a rough alignment image  $\mathbf{x}_a$

# Deep Generative Projection





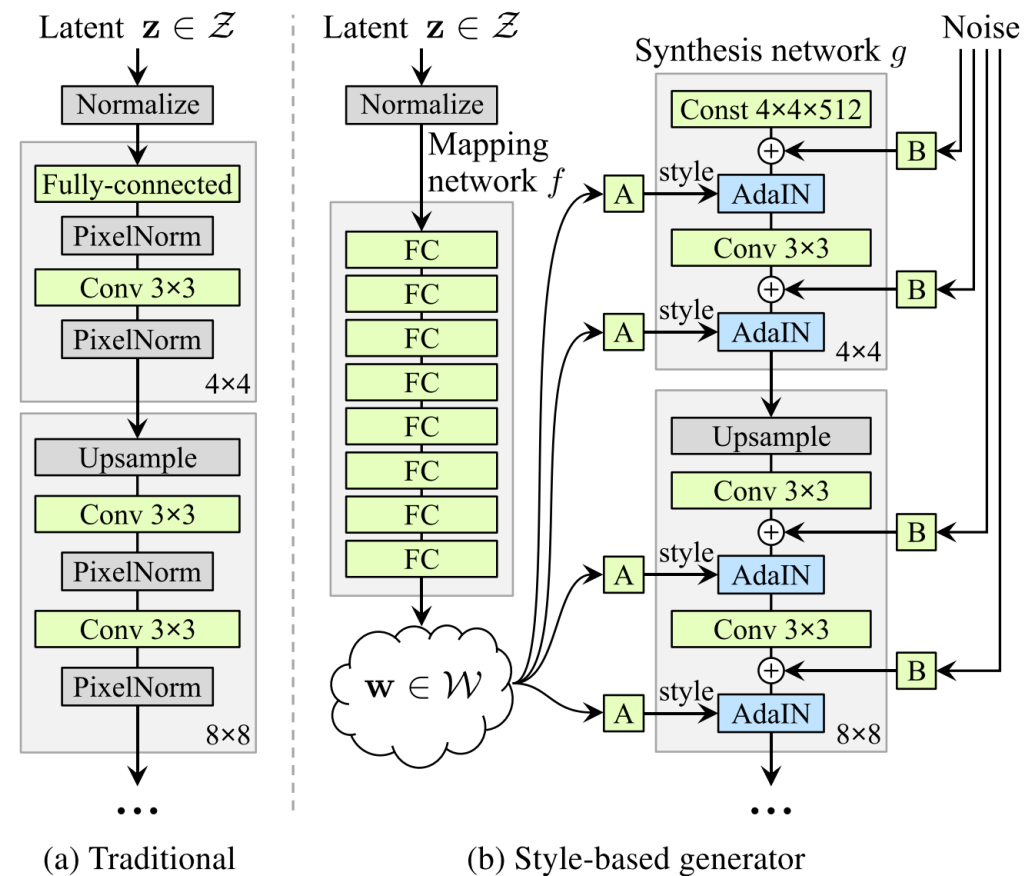
# Projection



- Rough alignment -> Image with high-level semantics (clothing category, model pose)
- Refine flaw of the rough alignment

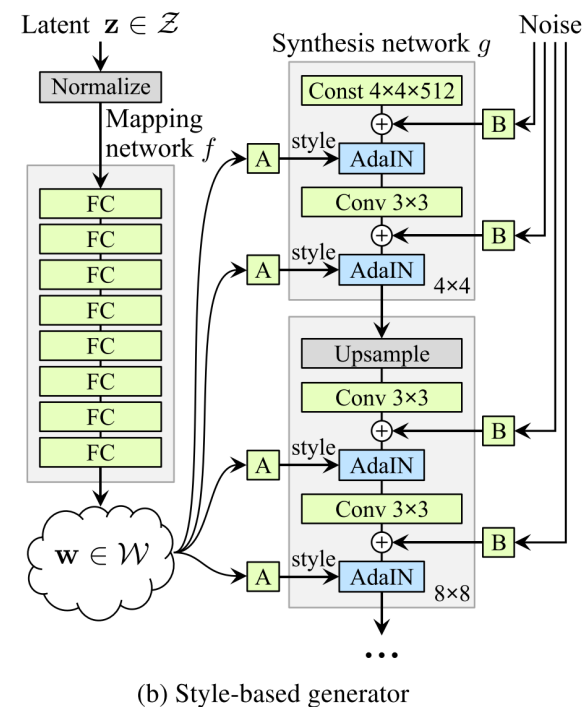
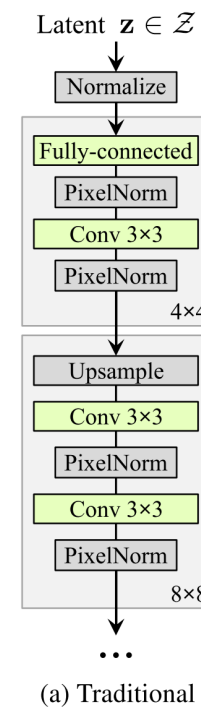
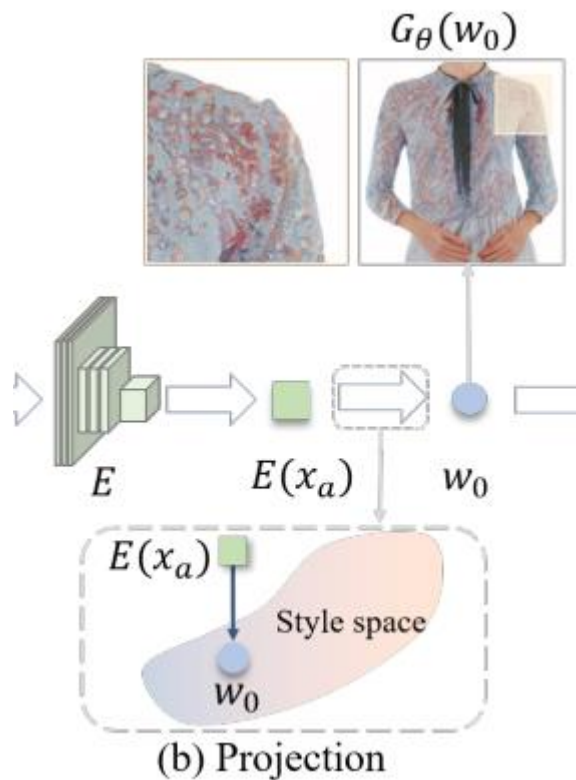
# Projection

- Train pre-trained StyleGAN
- E-Shop Fashion(ESF) dataset
  - 180,000 clothing model images
  - 512 x 512
- FID 2.16



# Projection

- Style space  $\mathcal{W}^+$
- Embedding images into high-density region of style space  
 -> To yield much more plausible synthesis



# Projection

- Encoder  $E$

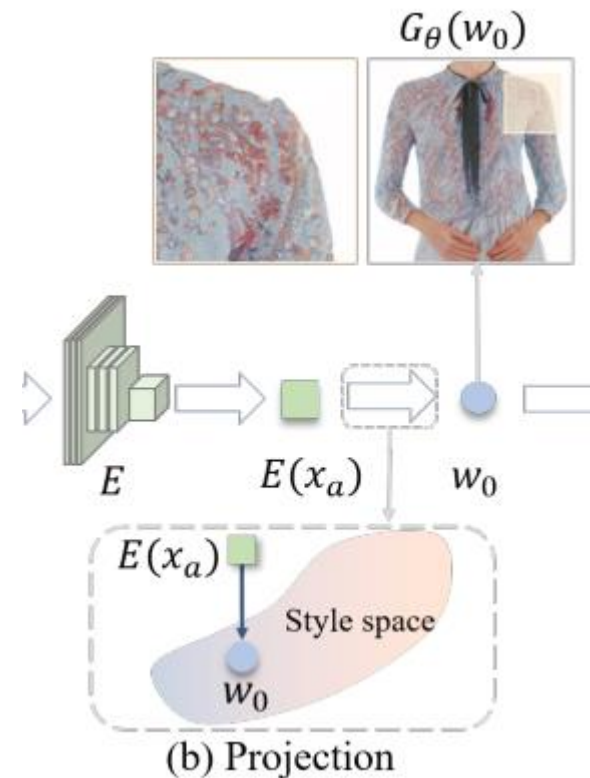
1. Compute PCA decomposition of 5M points on  $\mathcal{W}+$
2. Get the mean value, covariance matrix, set of principal components

$$\mu \quad \Sigma \quad Q = (q_1, \dots, q_n)$$

3. Predicts series of principal strengths  $s = (s_1, s_2, \dots, s_n)^T$

$$\begin{aligned} s &= E(x_a), \\ w_0 &= Tr(q_1 s_1 \sqrt{\sigma_1} + \dots q_n s_n \sqrt{\sigma_n}) + \mu \\ &= Q \Lambda^{\frac{1}{2}} Tr(s) + \mu, \end{aligned} \quad Tr(v) = \begin{cases} v, & \|v\|_2 < \psi, \\ \frac{v}{\|v\|_2} \psi, & \|v\|_2 \geq \psi. \end{cases}$$

$$w = P(x_a)$$



# Semantic search

- Optimization problem

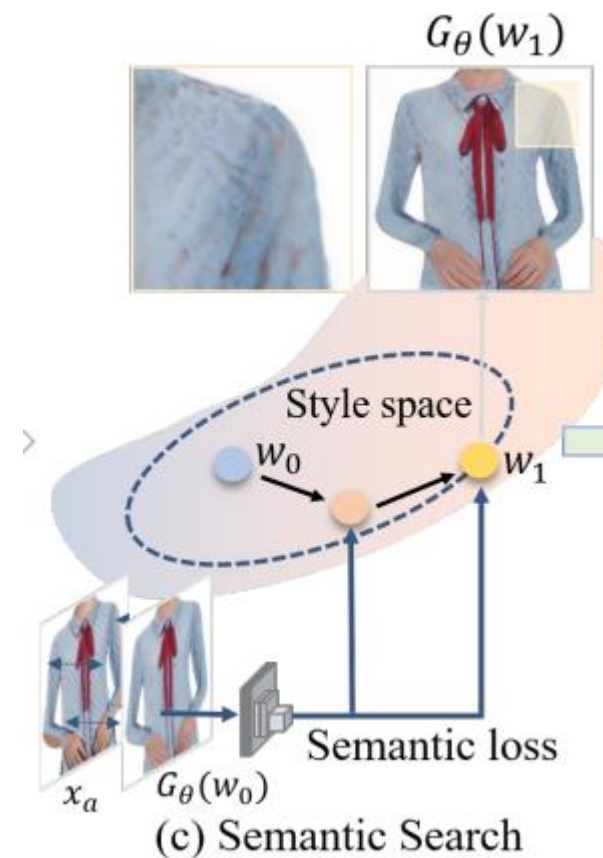
$$\min_{\mathbf{w} \in \mathcal{C}} \eta_p l_p + \eta_f l_f + \eta_{attr} l_{attr} + \eta_{adv} l_{adv},$$

$$l_p = \|W * \mathbf{G}(\mathbf{w}) - W * \mathbf{x}_a\|_2^2,$$

$$l_f = \|V(W * \mathbf{G}(\mathbf{w})) - V(W * \mathbf{x}_a)\|_2^2,$$

$$l_{attr} = \|R(W * \mathbf{G}(\mathbf{w})) - R(W * \mathbf{x}_a)\|_2^2,$$

$$l_{adv} = \log[1 - D(\mathbf{G}(\mathbf{w}))],$$





# Semantic search

- Optimization problem

$$\min_{\mathbf{w} \in \mathcal{C}} \eta_p l_p + \eta_f l_f + \eta_{attr} l_{attr} + \eta_{adv} l_{adv},$$

$$l_p = \|W * \mathbf{G}(\mathbf{w}) - W * \mathbf{x}_a\|_2^2,$$

$$l_f = \|\mathbf{V}(W * \mathbf{G}(\mathbf{w})) - \mathbf{V}(W * \mathbf{x}_a)\|_2^2,$$

$$l_{attr} = \|\mathbf{R}(W * \mathbf{G}(\mathbf{w})) - \mathbf{R}(W * \mathbf{x}_a)\|_2^2,$$

$$l_{adv} = \log[1 - D(\mathbf{G}(\mathbf{w}))],$$



Rough  
Alignment

Projection

$$W_{ij} = \begin{cases} 1 - \exp(-d((i, j), \partial I)^2), & I(ij) = 1, \\ 0, & I(ij) = 0. \end{cases}$$

# Semantic search

- Optimization problem

$$\begin{aligned} \min_{\mathbf{w} \in \mathcal{C}} \quad & \eta_p l_p + \eta_f l_f + \eta_{attr} l_{attr} + \eta_{adv} l_{adv}, \\ l_p = \quad & \|W * \mathbf{G}(\mathbf{w}) - W * \mathbf{x}_a\|_2^2, \\ l_f = \quad & \|\mathbf{V}(W * \mathbf{G}(\mathbf{w})) - \mathbf{V}(W * \mathbf{x}_a)\|_2^2, \\ l_{attr} = \quad & \|\mathbf{R}(W * \mathbf{G}(\mathbf{w})) - \mathbf{R}(W * \mathbf{x}_a)\|_2^2, \\ l_{adv} = \quad & \log[1 - D(\mathbf{G}(\mathbf{w}))], \end{aligned}$$

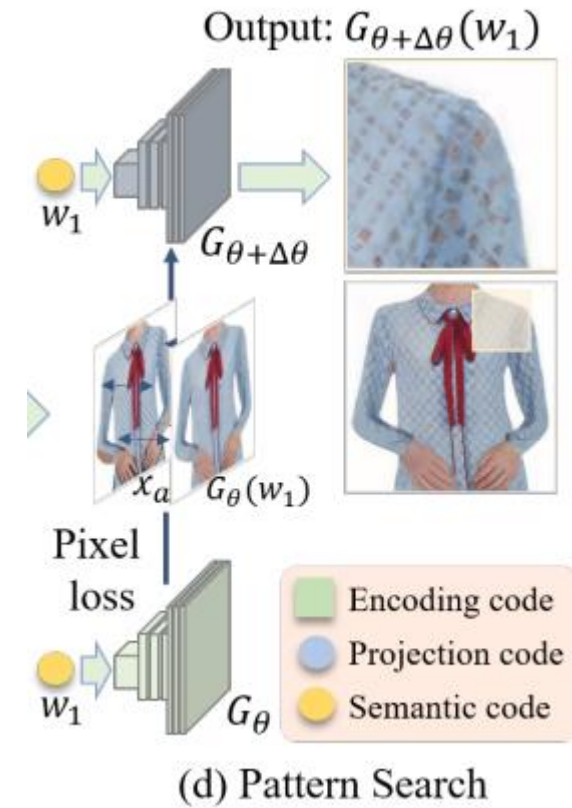
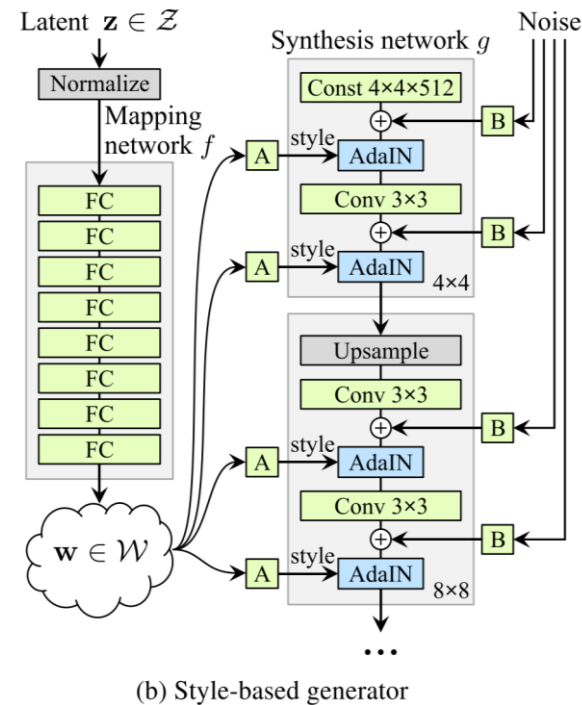
- Constrain  $\mathcal{C}$

$$\begin{aligned} \mathbf{w}_{k+1} &= \arg \min_{\mathbf{w} \in \mathcal{C}} \|\mathbf{w}_{k+1} - \mathbf{w}\| \\ &= \begin{cases} \mathbf{w}_0 + 4 \frac{\mathbf{w}_{k+1} - \mathbf{w}}{\|\mathbf{w}_{k+1} - \mathbf{w}\|_2}, & \|\mathbf{w}_{k+1} - \mathbf{w}_m\|_2 > 4, \\ \mathbf{w}_{k+1}, & \|\mathbf{w}_{k+1} - \mathbf{w}\|_2 \leq 4. \end{cases} \end{aligned}$$

# Pattern search

- Optimization problem

$$\min_{\theta \in B(\theta_0, 4)} \eta_p \|W * G_\theta(w) - W * x_a\|_2 + \log(1 - D(G_\theta(w))),$$



# Experiments

Table 1. Numerical metrics of DGP, ACGPN, PF-AFN, and VITON-HD on CMI and MPV datasets. ↓ indicates lower is better.

Methods	CMI		MPV	
	FID↓	SWD↓	FID↓	SWD↓
ACGPN	137.9	121.3	81.1	90.4
PF-AFN	97.3	76.7	67.8	67.1
VITON-HD	87.5	56.1	<b>40.6</b>	52.7
DGP (Ours)	<b>51.6</b>	<b>22.4</b>	48.4	<b>36.7</b>

# Experiments



Figure 8. Comparison on the CMI and MPV datasets. The supervised competitor methods are basically less appealing, and perform especially poorly on complicated clothing like coats.



# Experiments

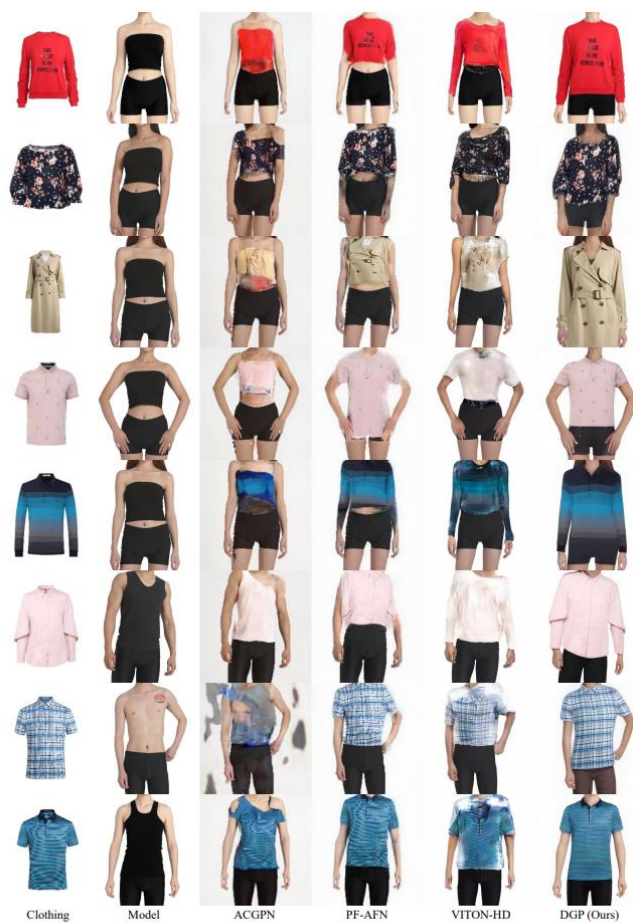


Figure S10. More visual results of qualitative comparison on the CMI dataset.