



南京航空航天大學
NANJING UNIVERSITY OF AERONAUTICS AND ASTRONAUTICS



Person Image Synthesis via Denoising Diffusion Model

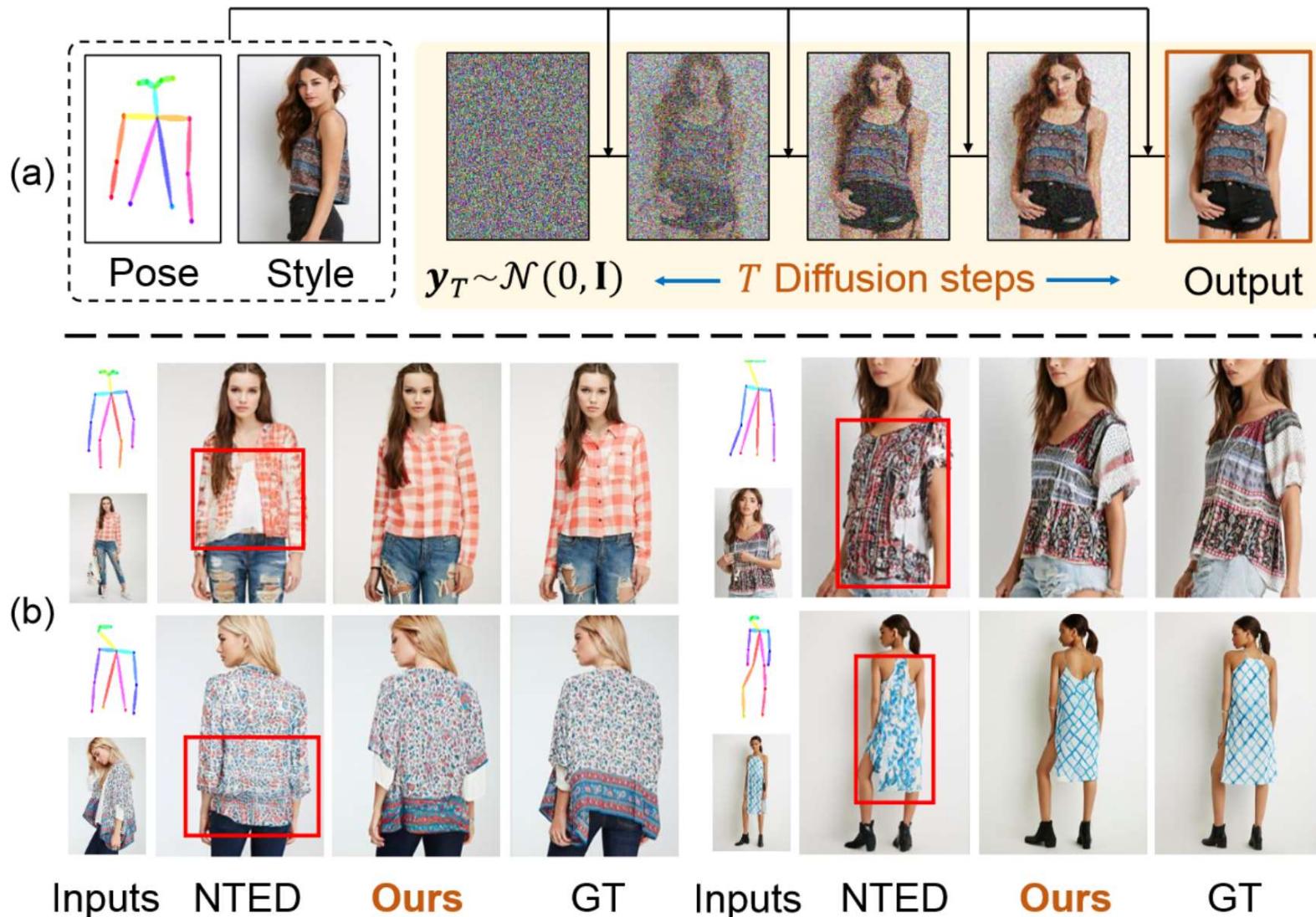
Ankan Kumar Bhunia¹ Salman Khan^{1,2} Hisham Cholakkal¹ Rao Muhammad Anwer^{1,4}
Jorma Laaksonen⁴ Mubarak Shah⁵ Fahad Shahbaz Khan^{1,3}

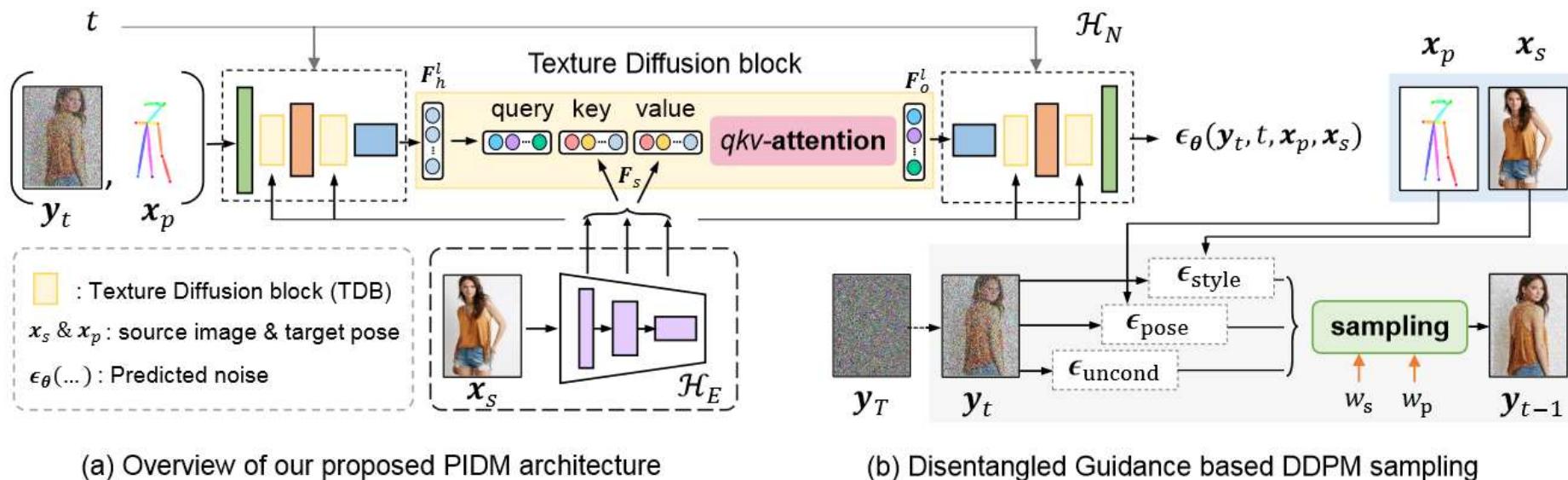
¹Mohamed bin Zayed University of AI, UAE ²Australian National University, Australia

³Linköping University, Sweden ⁴Aalto University, Finland ⁵University of Central Florida, USA

汇报人：李雅超

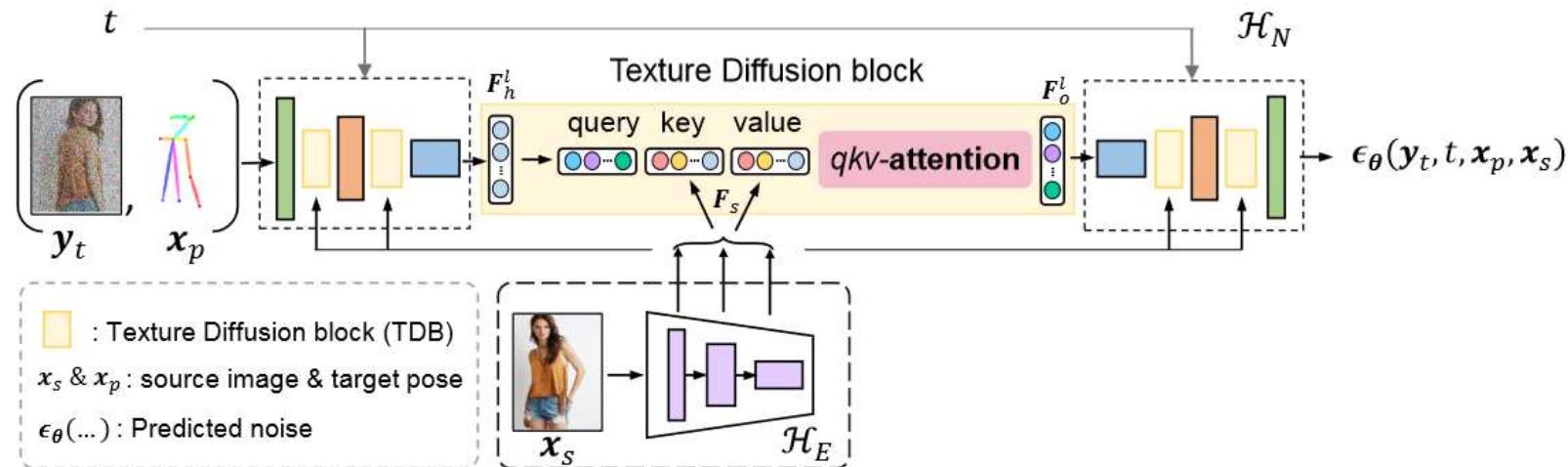
Bhunia, Ankan Kumar, et al. "Person image synthesis via denoising diffusion model." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2023.







Texture-Conditioned Diffusion Model



前向扩散分布:

$$q(y_t|y_{t-1}) = \mathcal{N}(y_t; \sqrt{1 - \beta_t}y_{t-1}, \beta_t\mathbf{I}). \quad (1)$$

逆向去噪分布:

$$\begin{aligned} p_\theta(y_{t-1}|y_t, x_p, x_s) &= \mathcal{N}(y_{t-1}; \mu_\theta(y_t, t, x_p, x_s), \\ &\quad \Sigma_\theta(y_t, t, x_p, x_s)). \end{aligned} \quad (2)$$

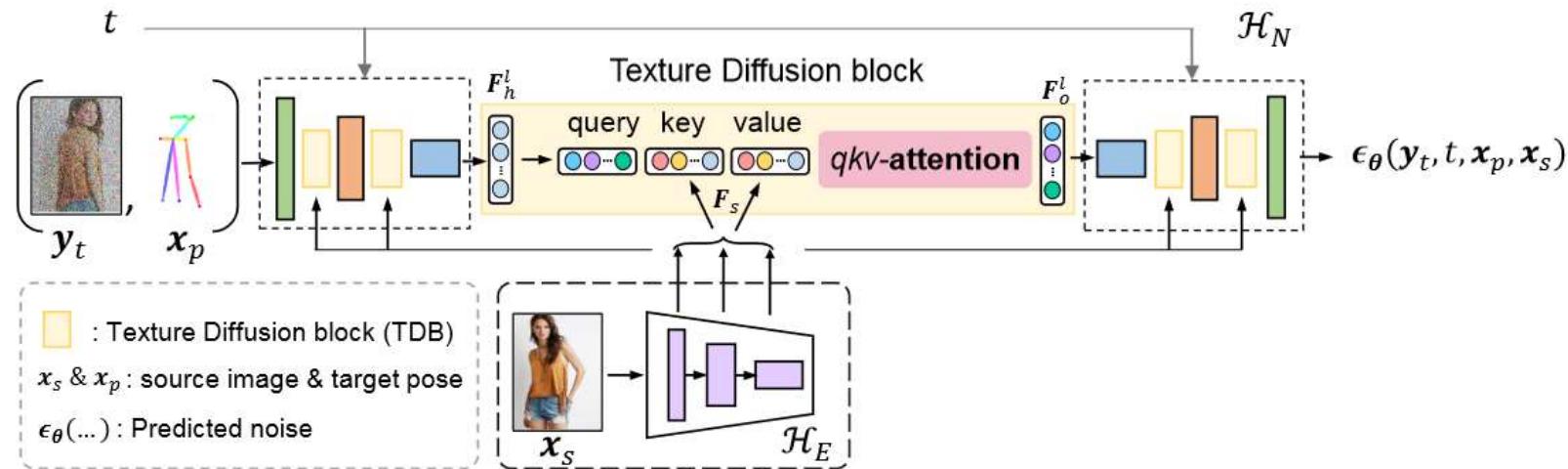
损失函数:

$$L_{\text{mse}} = \mathbb{E}_{t \sim [1, T], y_0 \sim q(y_0), \epsilon} \|\epsilon - \epsilon_\theta(y_t, t, x_p, x_s)\|^2. \quad (3)$$

$$L_{\text{hybrid}} = L_{\text{mse}} + L_{\text{vib}}. \quad (4)$$



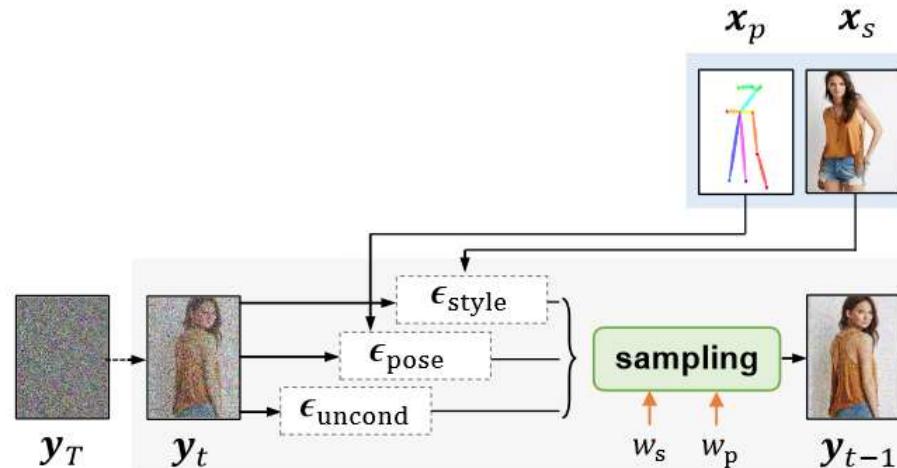
Texture Diffusion block



$$\begin{aligned} \mathbf{Q} &= \phi_q^l(\mathbf{F}_h^l), \quad \mathbf{K} = \phi_k^l(\mathbf{F}_s), \quad \mathbf{V} = \phi_v^l(\mathbf{F}_s) \\ \mathbf{F}_{att}^l &= \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{C}}, \quad \mathbf{F}_o^l = \mathbf{W}^l \text{softmax}(\mathbf{F}_{att}^l)\mathbf{V} + \mathbf{F}_h^l, \quad (5) \end{aligned}$$



Disentangled Guidance based Sampling



$$\epsilon_{\text{uncond}} = \epsilon_{\theta}(y_t, t, \emptyset, \emptyset)$$

$$\epsilon_{\text{pose}} = \epsilon_{\theta}(y_t, t, x_p, \emptyset) - \epsilon_{\text{uncond}}$$

$$\epsilon_{\text{style}} = \epsilon_{\theta}(y_t, t, \emptyset, x_s) - \epsilon_{\text{uncond}}$$

解耦的条件噪音：

$$\epsilon_{\text{cond}} = \epsilon_{\text{uncond}} + w_p \epsilon_{\text{pose}} + w_s \epsilon_{\text{style}}, \quad (6)$$

原始的无分类引导下的条件噪音：

$$\begin{aligned}\bar{\epsilon}_{\theta}(\mathbf{x}_t, t, y) &= \epsilon_{\theta}(\mathbf{x}_t, t, y) - \sqrt{1 - \bar{\alpha}_t} w \nabla_{\mathbf{x}_t} \log p(y | \mathbf{x}_t) \\ &= \epsilon_{\theta}(\mathbf{x}_t, t, y) + w (\epsilon_{\theta}(\mathbf{x}_t, t, y) - \epsilon_{\theta}(\mathbf{x}_t, t)) \\ &= (w + 1) \epsilon_{\theta}(\mathbf{x}_t, t, y) - w \epsilon_{\theta}(\mathbf{x}_t, t)\end{aligned}$$

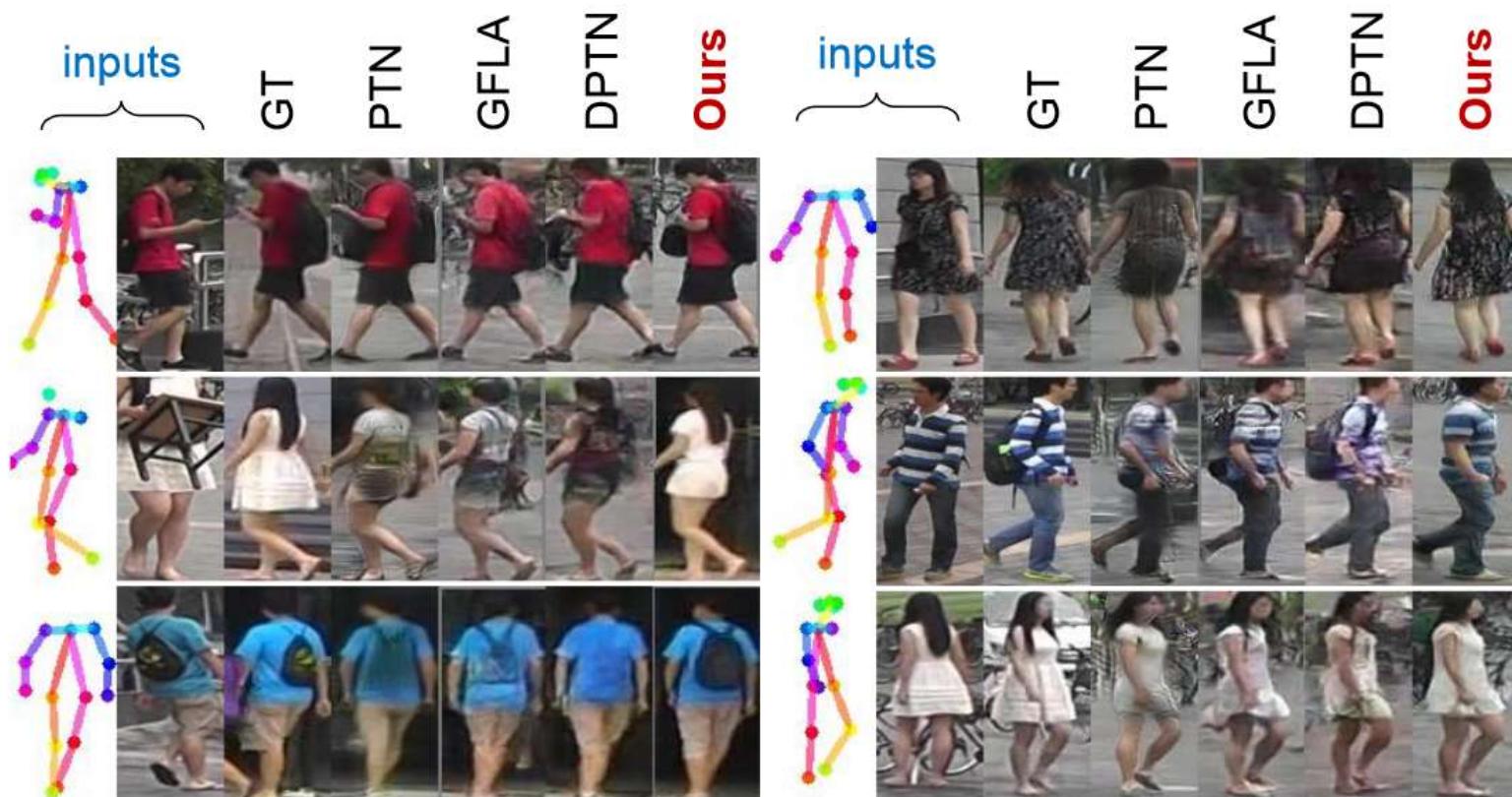


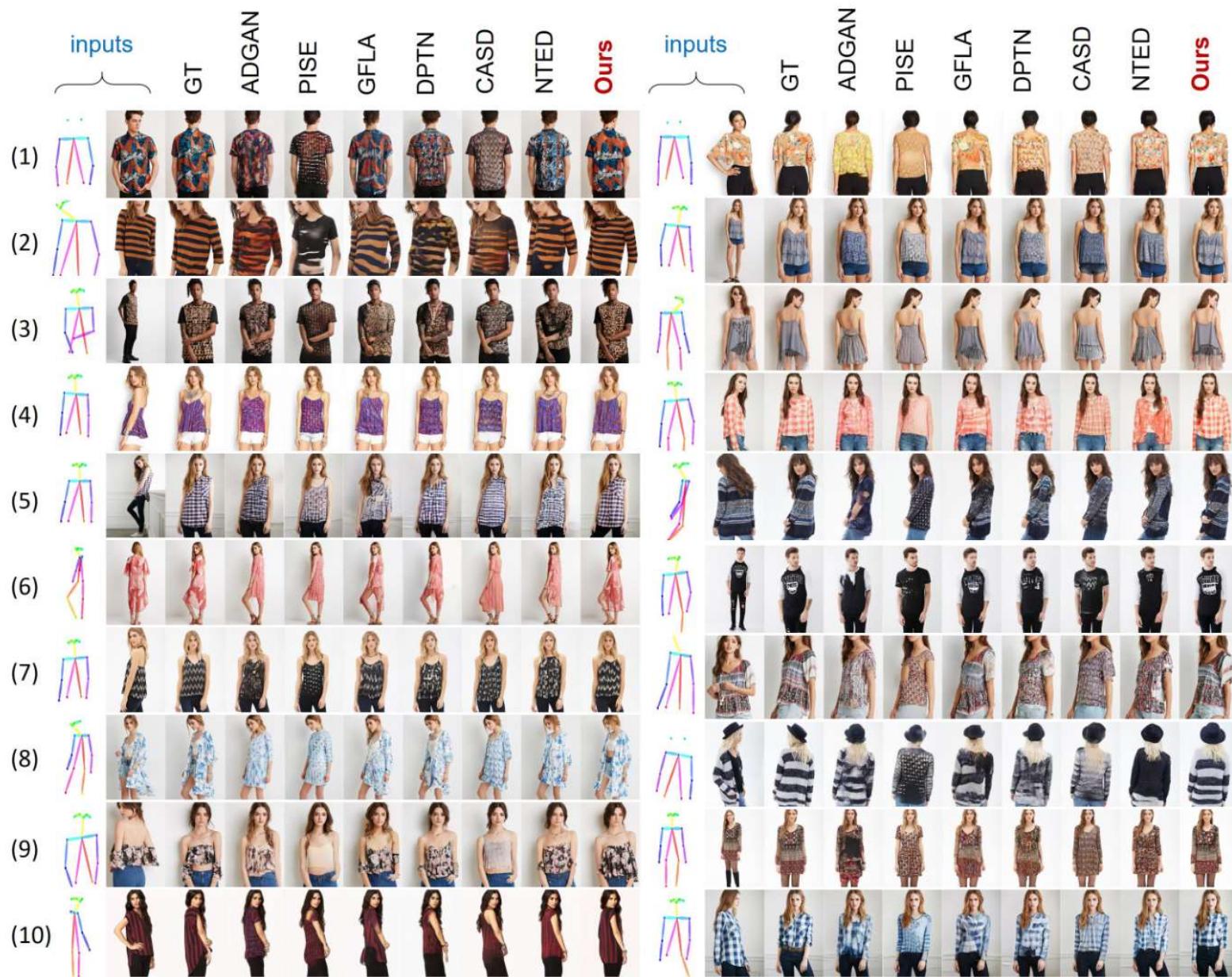
定量实验结果

Dataset	Methods	FID(\downarrow)	SSIM(\uparrow)	LPIPS(\downarrow)
DeepFashion [11] (256 × 176)	Def-GAN [20]	18.457	0.6786	0.2330
	PATN [30]	20.751	0.6709	0.2562
	ADGAN [14]	14.458	0.6721	0.2283
	PISE [23]	13.610	0.6629	0.2059
	GFLA [19]	10.573	0.7074	0.2341
	DPTN [24]	11.387	0.7112	0.1931
	CASD [28]	11.373	0.7248	0.1936
	NTED [18]	8.6838	0.7182	0.1752
	PIDM (Ours)	6.3671	0.7312	0.1678
DeepFashion [11] (512 × 352)	CocosNet2 [29]	13.325	0.7236	0.2265
	NTED [18]	7.7821	0.7376	0.1980
	PIDM (Ours)	5.8365	0.7419	0.1768
Market-1501 [27] (128 × 64)	Def-GAN [20]	25.364	0.2683	0.2994
	PTN [30]	22.657	0.2821	0.3196
	GFLA [19]	19.751	0.2883	0.2817
	DPTN [24]	18.995	0.2854	0.2711
	PIDM (Ours)	14.451	0.3054	0.2415



可视化对比结果



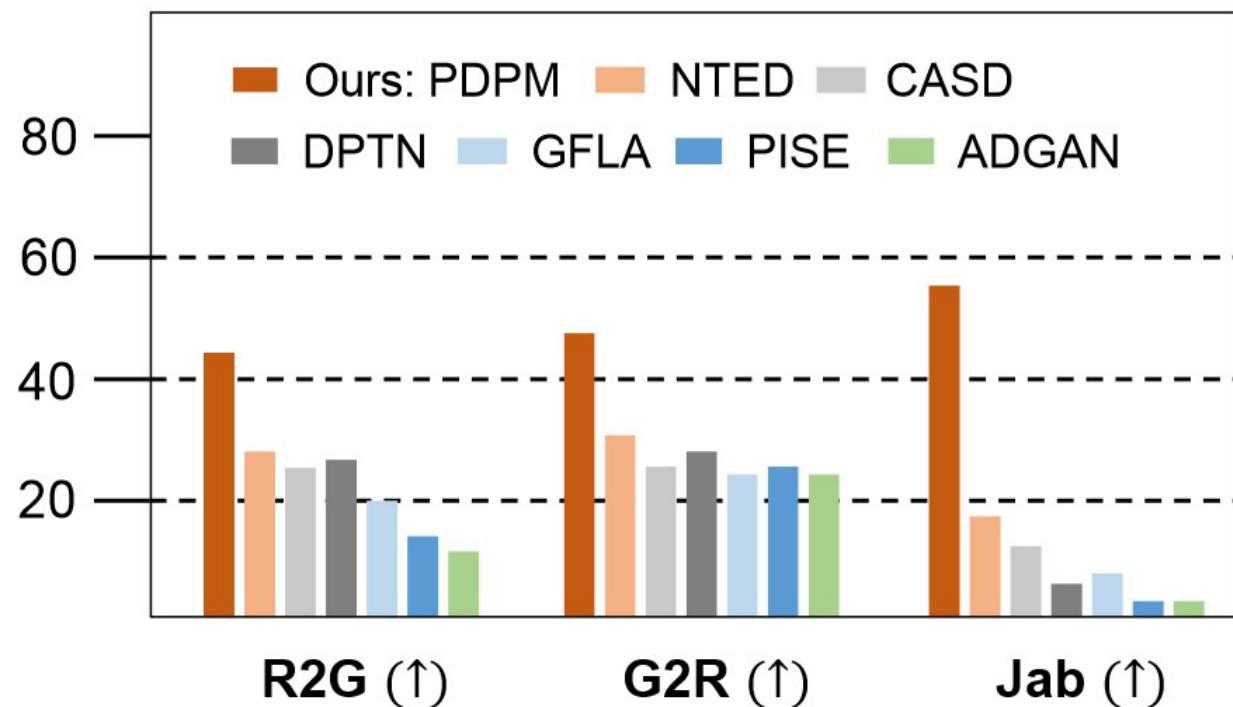




User Study结果

30 generated images

30 real-world images



R2G: real-world to generate

G2R: generate to real-world

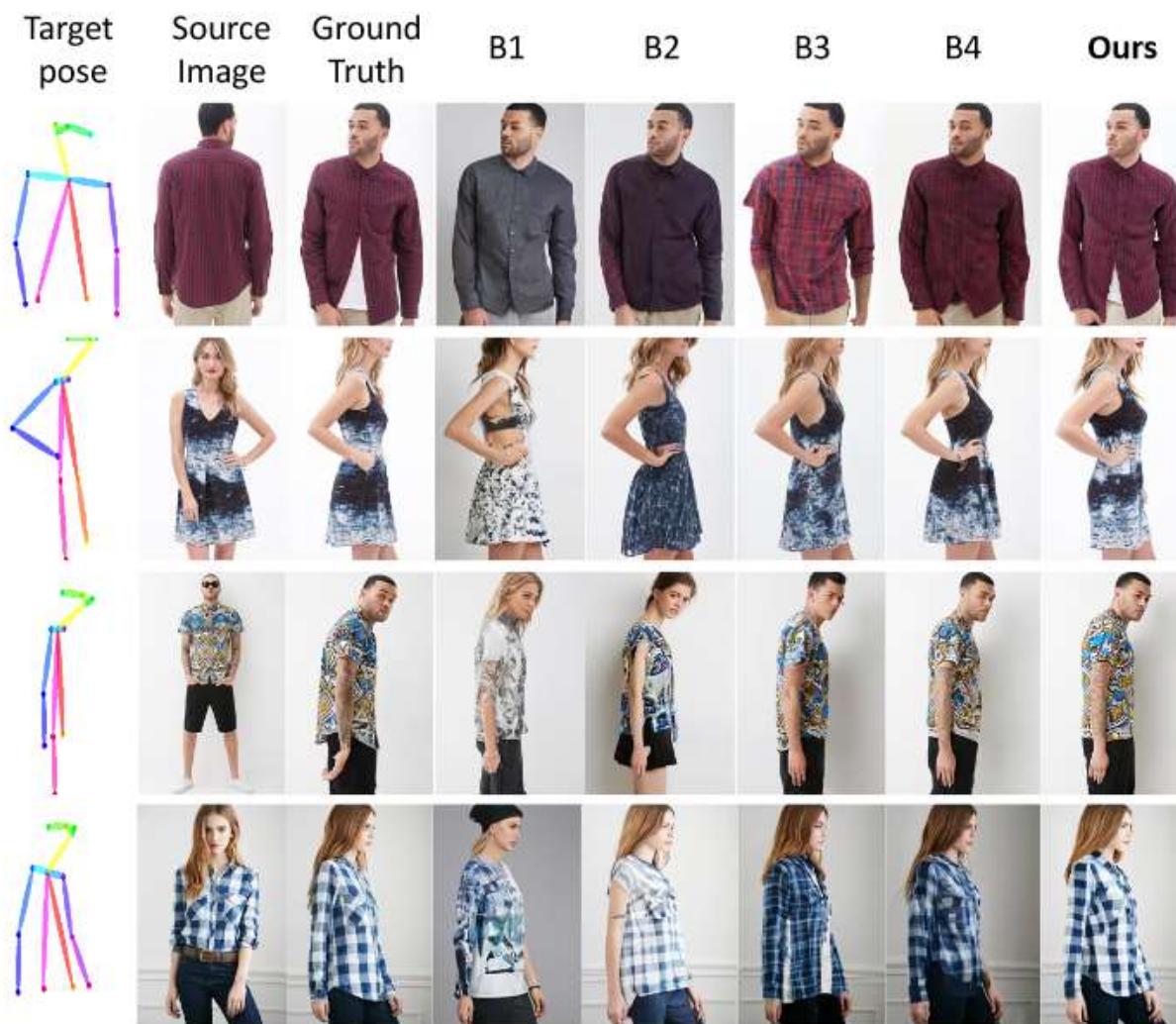
Jab: the percentage of images considered the best among all models



消融实验

Methods	FID(\downarrow)	SSIM(\uparrow)	LPIPS(\downarrow)
B1: Baseline [†] (concat)	10.813	0.6911	0.2112
B2: Baseline	9.8510	0.7005	0.1983
B3: Baseline + TDB	7.5133	0.7178	0.1870
B4: Baseline + TDB + CF-guidance	6.8176	0.7195	0.1769
Ours: Baseline + TDB + DCF-guidance	6.3671	0.7312	0.1676

消融实验



Appearance Control and Editing

