

# Multistage Adversarial Losses for Pose-Based Human Image Synthesis

CVPR 2018

Center for Research on Intelligent Perception and Computing (CRIPAC),  
National Laboratory of Pattern Recognition (NLPR)  
Center for Excellence in Brain Science and Intelligence Technology  
(CEBSIT), Institute of Automation, Chinese Academy of Sciences (CASIA)  
University of Chinese Academy of Sciences (UCAS)

# Motivation



Input

cGANs[I16]

VSA[25]

PG<sup>2</sup>[I13]

Ours

GT



# Method Overview

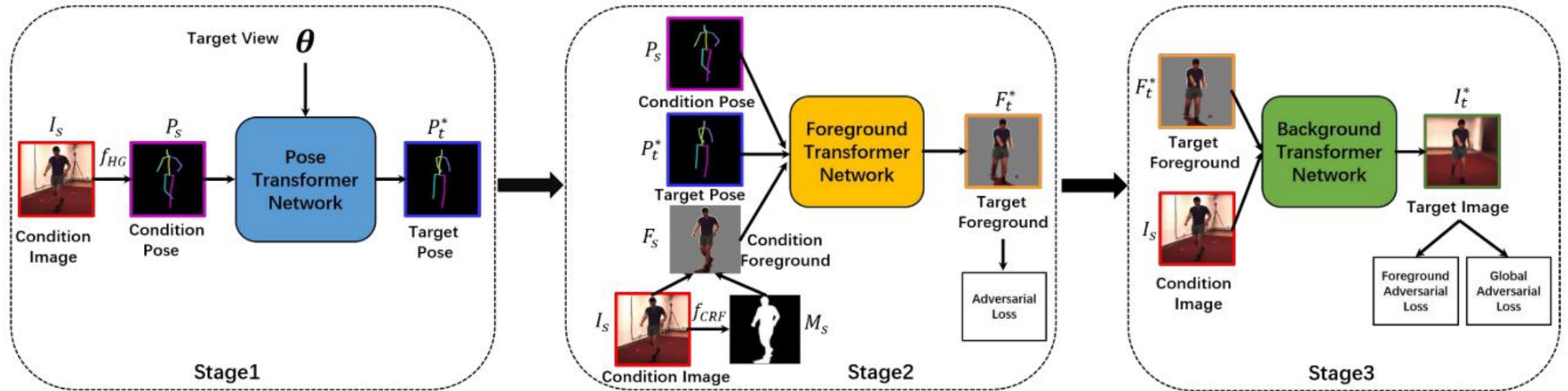


Figure 2. The overall pipeline of our multistage approach which contains three transformer networks for three stages. In the first stage, the pose transformer network synthesizes a novel view 2D pose. Then, the foreground transformer network synthesizes the target foreground image in the second stage. Finally, the background transformer network generates the target image.  $f_{HG}$  and  $f_{CRF}$  donate the stacked hourglass networks [18] and the CRF-RNN [30] for pose estimation from image and foreground segmentation, respectively.

# Pose Transformer Network

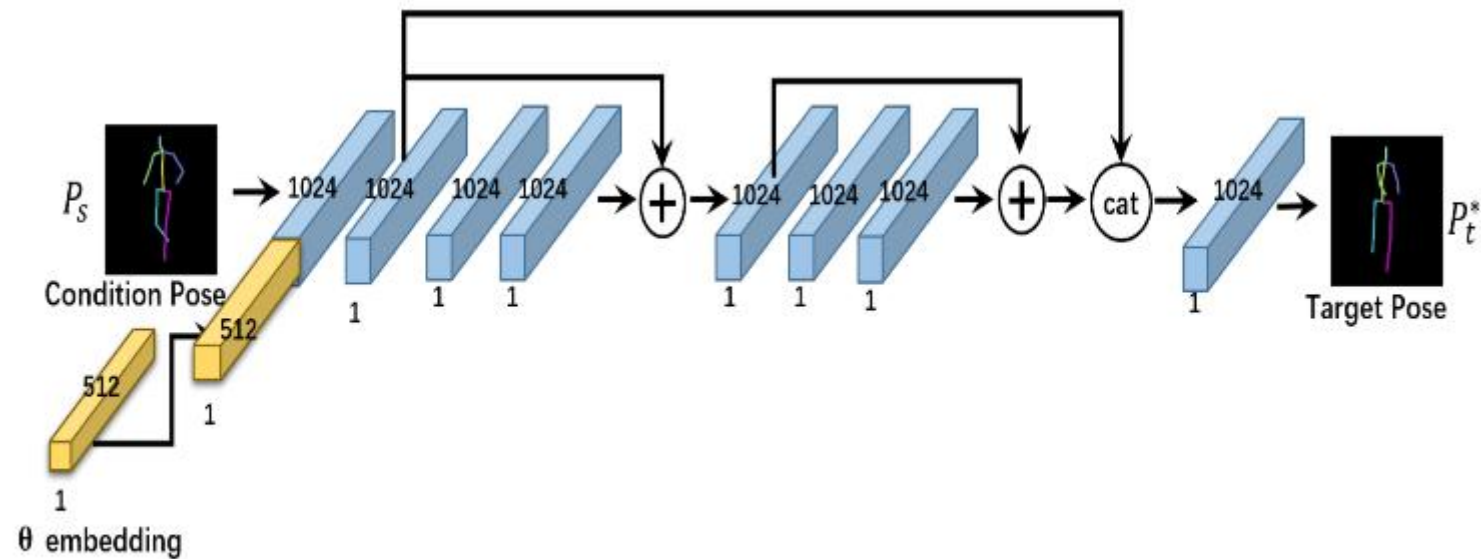


Figure 3. The architecture of the pose transformer network.

$$P_t^* = G_p(P_s, \theta)$$

$$\mathcal{L}^1 = \sum_i^N \left\| P_t^{*i} - P_t^i \right\|_2^2$$

# Foreground Transformer Network

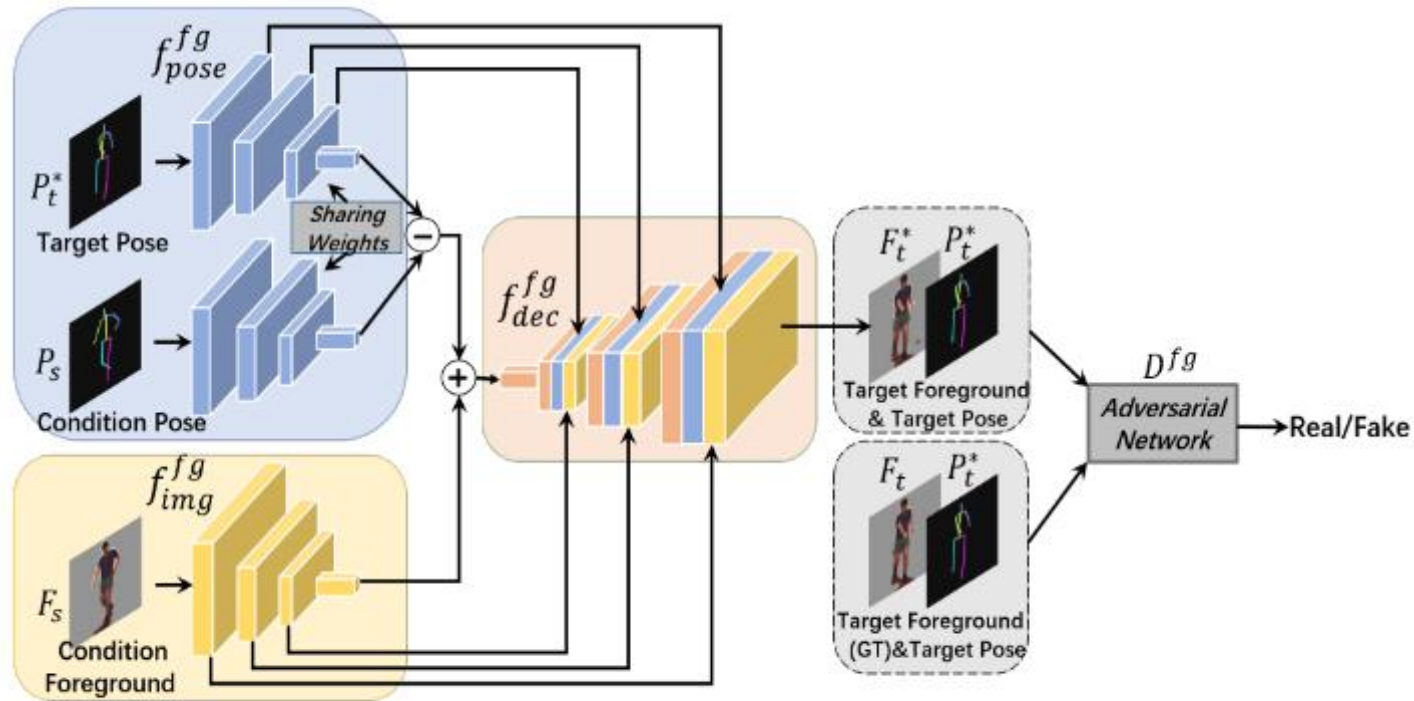


Figure 4. The architecture of the foreground transformer network.

$$F_t^* = f_{dec}^{fg}(f_{pose}^{fg}(P_t^*) - f_{pose}^{fg}(P_s) + f_{img}^{fg}(F_s))$$

# Foreground Transformer Networks

$$\mathcal{L}^2 = \alpha_f \mathcal{L}_{fg}^2 + \beta_f \mathcal{L}_{bg}^2 + \mathcal{L}_{gen}^2$$

$$\begin{aligned} \mathcal{L}_{fg}^2 &= \|F_t \odot M_t - F_t^* \odot M_t\|_1 \\ &= \frac{1}{\sum_{M_t^{i,j}=1} M_t^{i,j}} \sum_{i,j} \left| (F_t^{i,j} - F_t^{*,i,j}) \times M_t^{i,j} \right| \end{aligned}$$

$$\mathcal{L}_{gen}^2 = -\log(D^{fg}([F_t^*, P_t^*]))$$

$$\begin{aligned} \mathcal{L}_{bg}^2 &= \|F_t \odot (1 - M_t) - F_t^* \odot (1 - M_t)\|_1 \\ &= \frac{1}{\sum_{M_t^{i,j}=0} (1 - M_t^{i,j})} \sum_{i,j} \left| (F_t^{i,j} - F_t^{*,i,j}) \times (1 - M_t^{i,j}) \right| \end{aligned}$$

$$\begin{aligned} \mathcal{L}_D^2 &= -\log(D^{fg}([F_t, P_t^*])) \\ &\quad - \log(1 - D^{fg}([F_t^*, P_t^*])) \end{aligned}$$



# Background transformer network

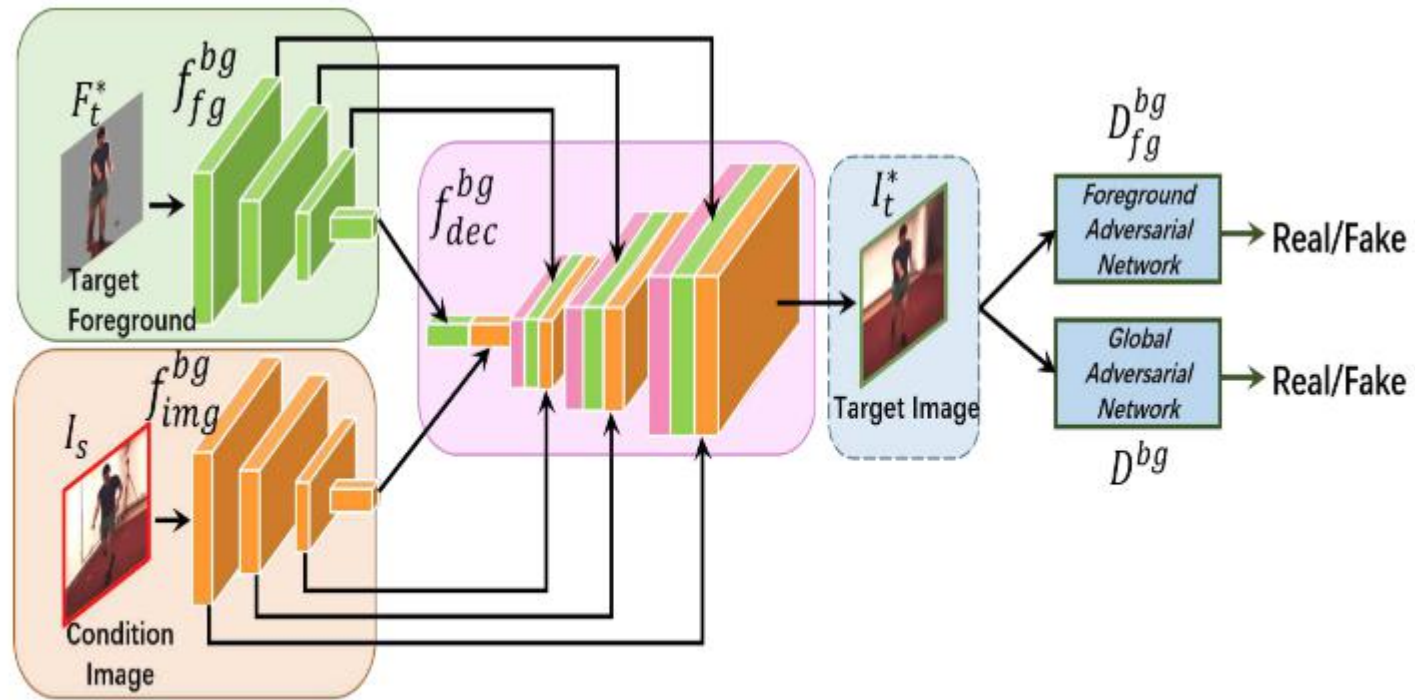


Figure 5. The architecture of the background transformer network.

# Background Transformer Networks

$$\mathcal{L}^3 = \alpha_b \mathcal{L}_{fg}^3 + \beta_b \mathcal{L}_{bg}^3 + \mathcal{L}_{gen_{fg}}^3 + \mathcal{L}_{gen}^3$$

$$\mathcal{L}_{gen_{fg}}^3 = -\log(D_{fg}^{bg}([I_t^* \odot M_t, P_t^*]))$$

$$\mathcal{L}_{gen}^3 = -\log(D^{bg}(I_t^*))$$

$$\mathcal{L}_{D_{fg}}^3 = -\log(D_{fg}^{bg}([I_t \odot M_t, P_t^*]))$$

$$-\log(1 - D_{fg}^{bg}([I_t^* \odot M_t, P_t^*]))$$

$$\mathcal{L}_D^3 = -\log(D^{bg}(I_t)) - \log(1 - D^{bg}(I_t^*))$$



# Experiment



Figure 6. Visualization of the synthesized images from three state-of-the-art methods, two baselines and our model. Our method achieves the best results with clear foreground and background.

# Experiment

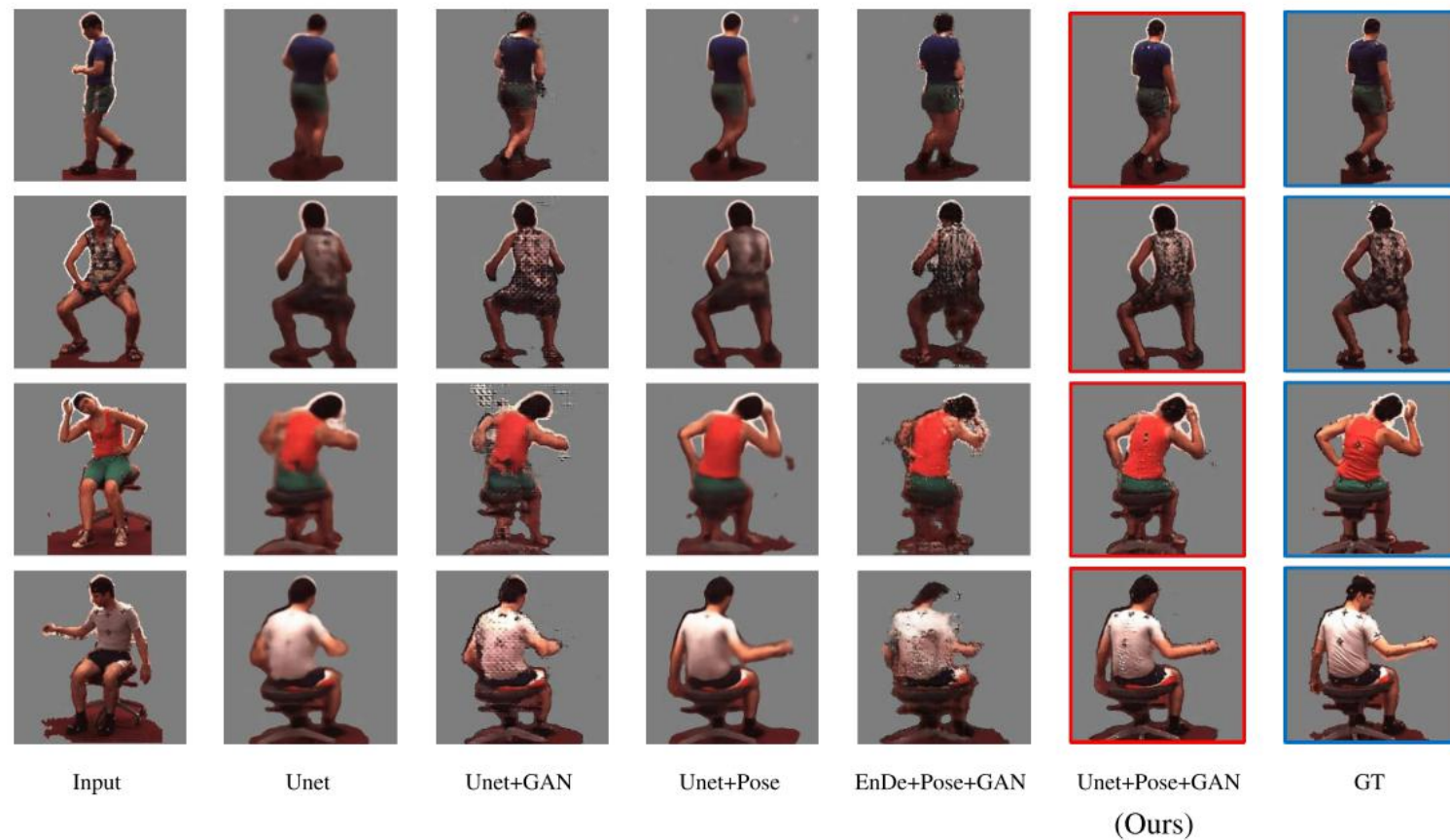
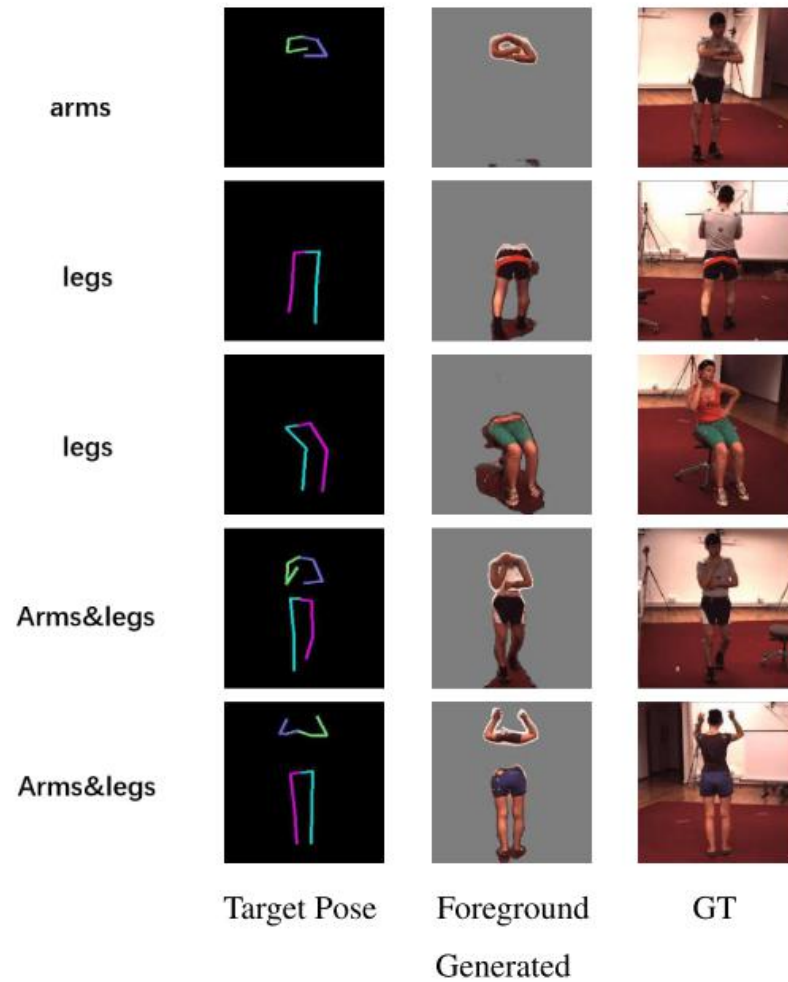


Figure 7. Visualization of the synthesized images from four foreground baselines and our foreground transformer network.

# Experiment





# Experiment

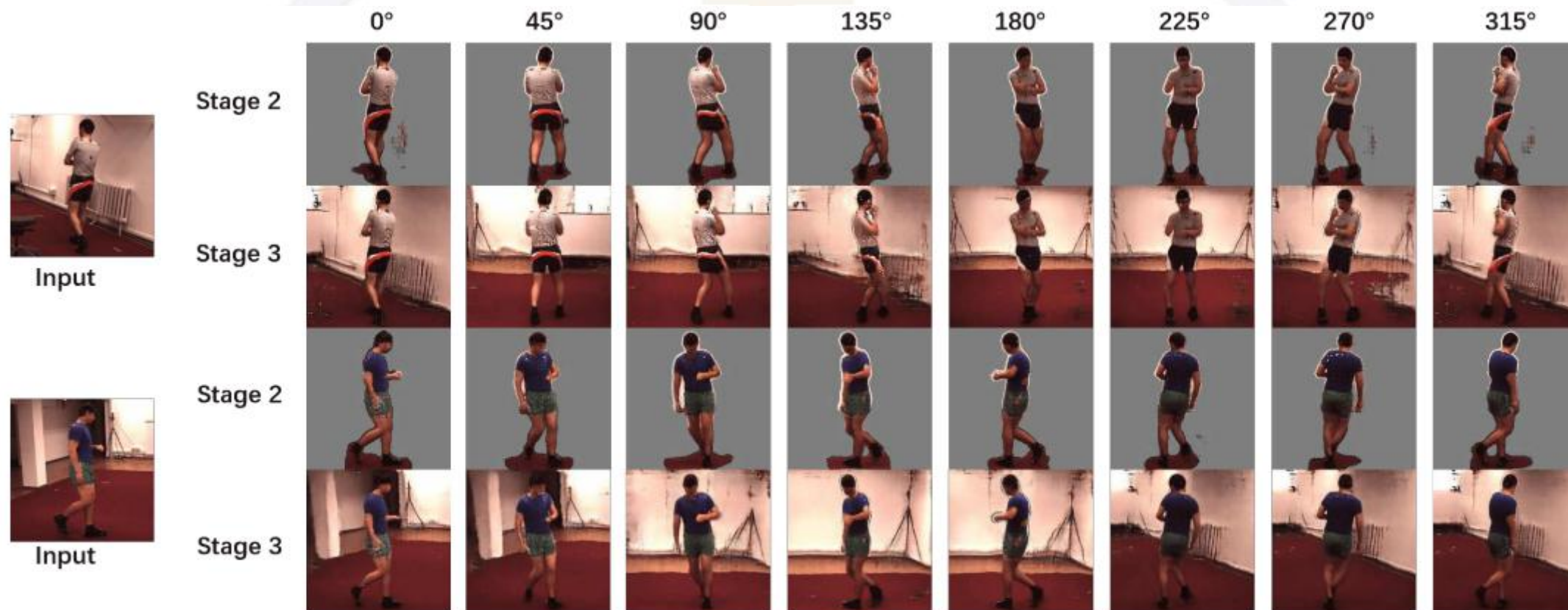


Figure 9. Multiview human images generated by our model.

# Experiment

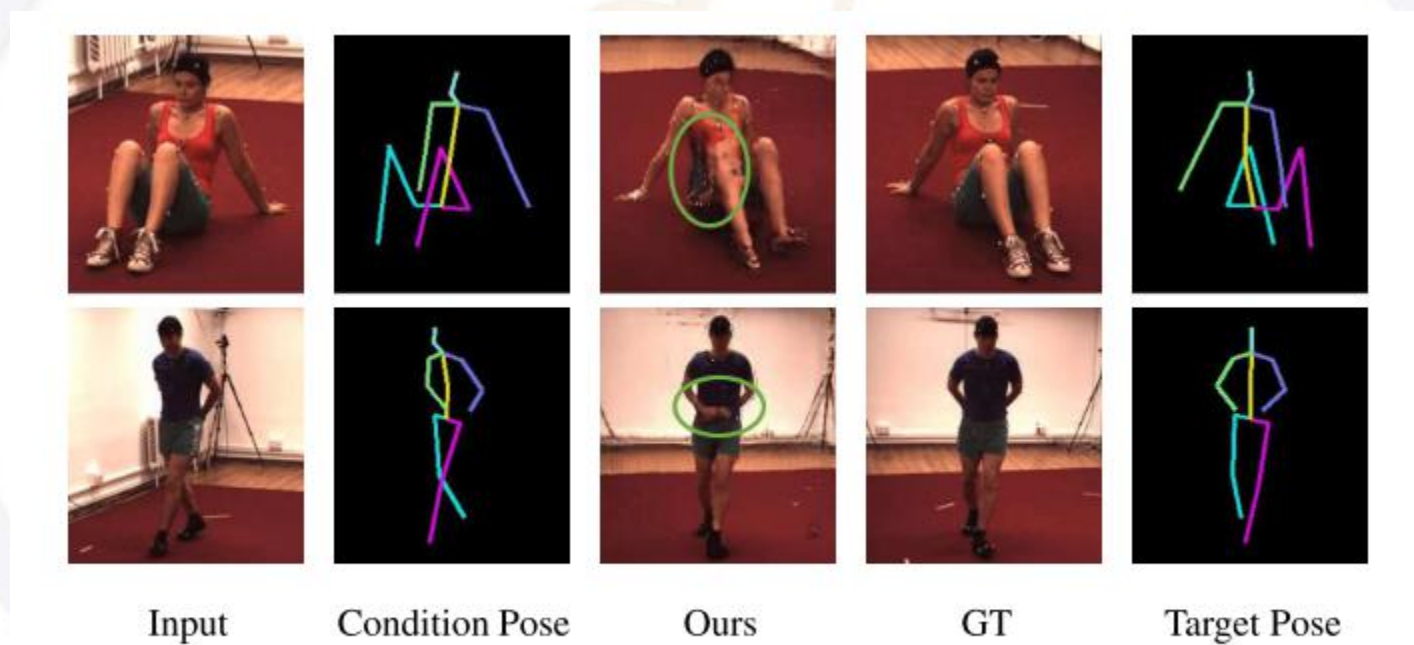


Figure 10. Two failure cases of our model.



THANKS