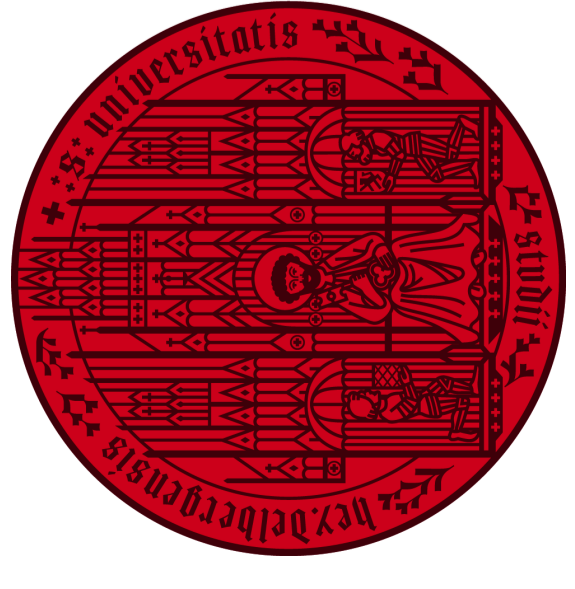


A Variational U-Net for Conditional Appearance and Shape Generation



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Heidelberg Collaboratory



for Image Processing

Patrick Esser* Ekaterina Sutter* Björn Ommer

{firstname.lastname}@iwr.uni-heidelberg.de

HCI, IWR, Heidelberg University

*equal contribution

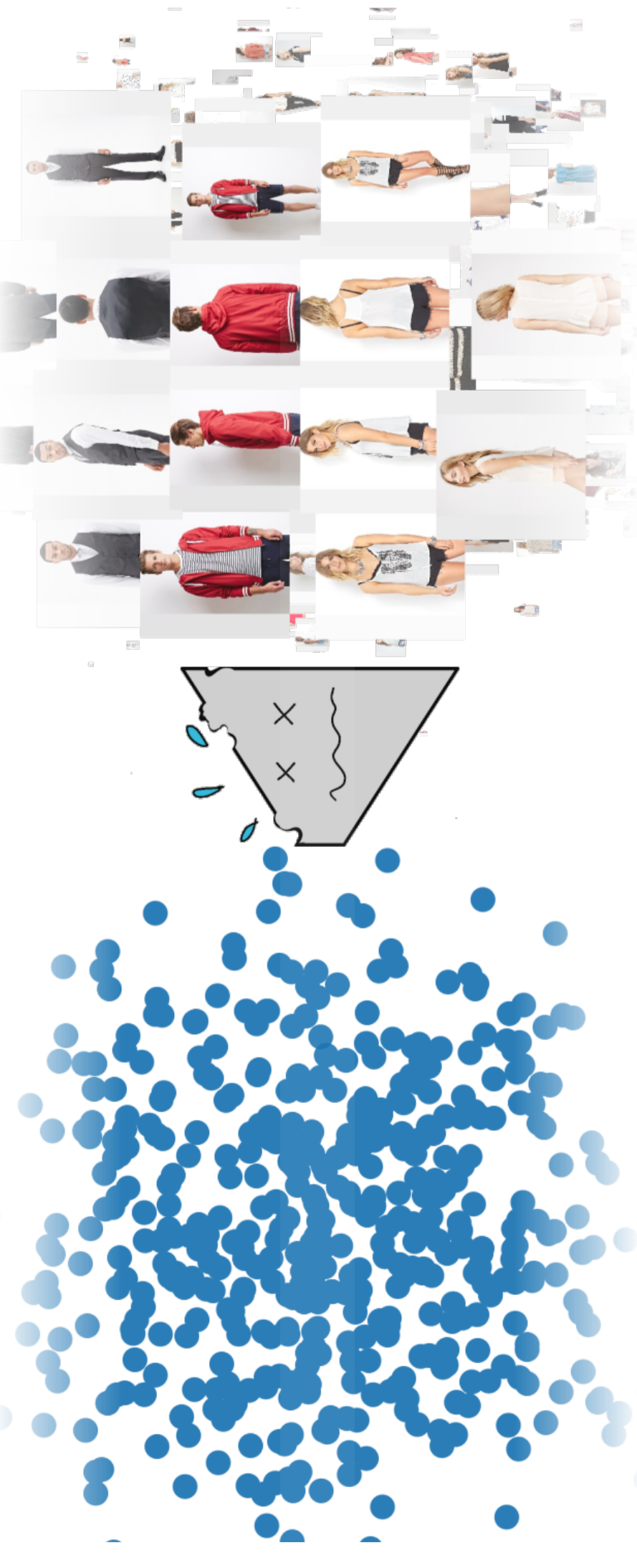


Interdisciplinary Center
for Scientific Computing



Motivation

Many previous methods try to generate all image variations in the same way.



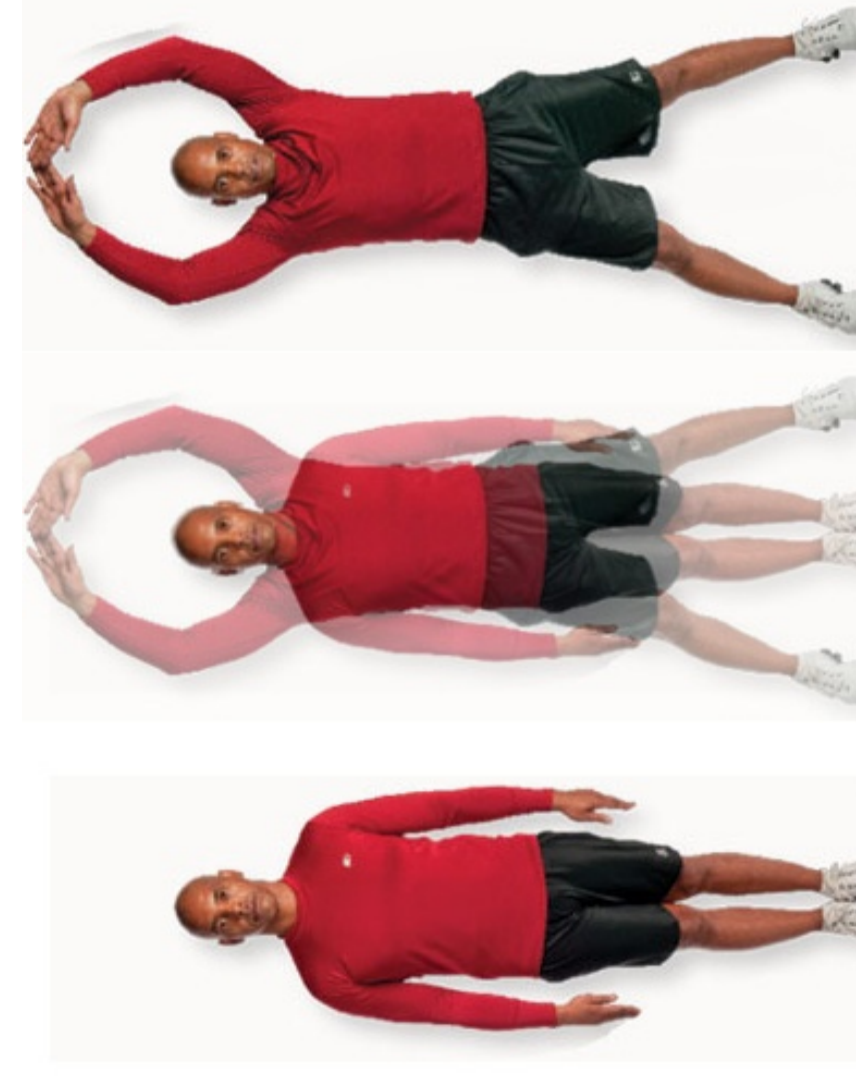
However, we can identify two groups of variations with different characteristics:

Variations in \rightarrow Appearance
color, texture, identity, ...



a linear interpolation between these changes works well

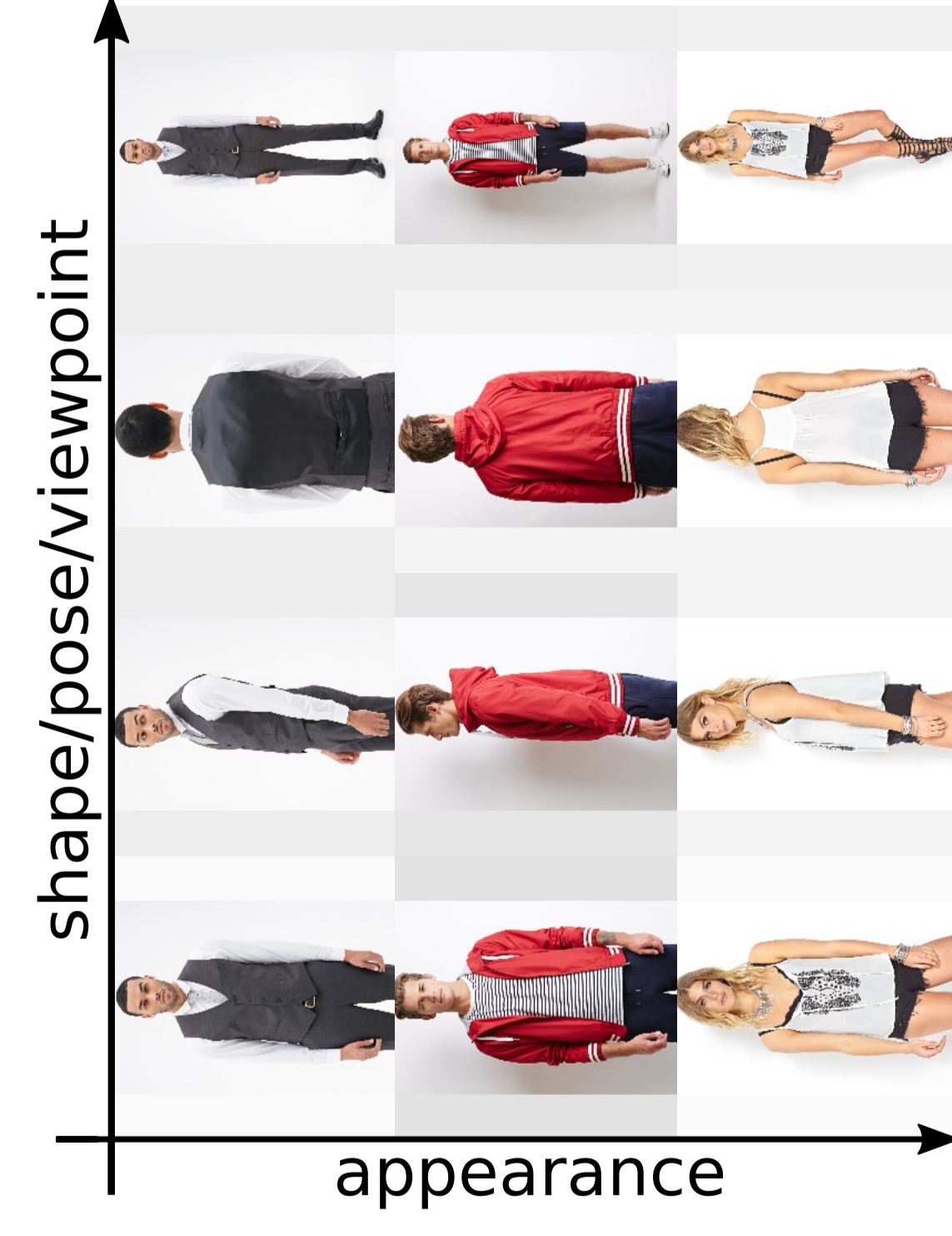
\rightarrow Shape
shape, pose, viewpoint, ...



deformations cannot be explained in the same way

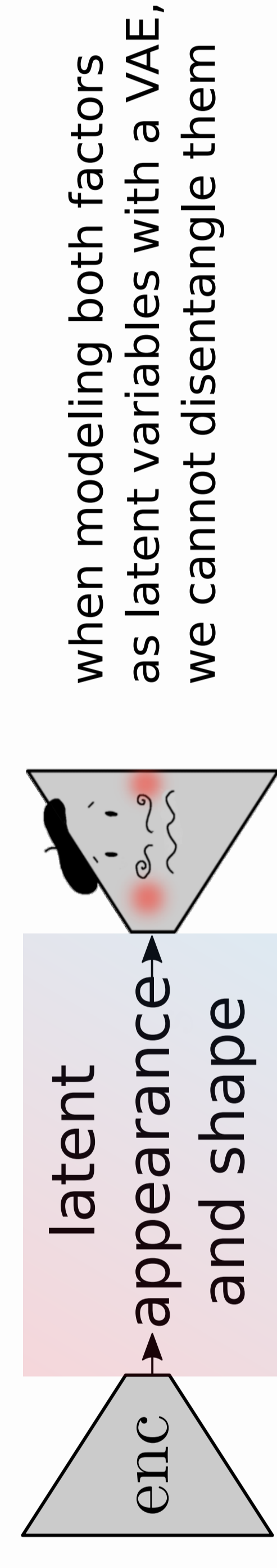
A model which understands this difference must be able to control appearance and shape separately.

How can we disentangle & generate appearance and shape?



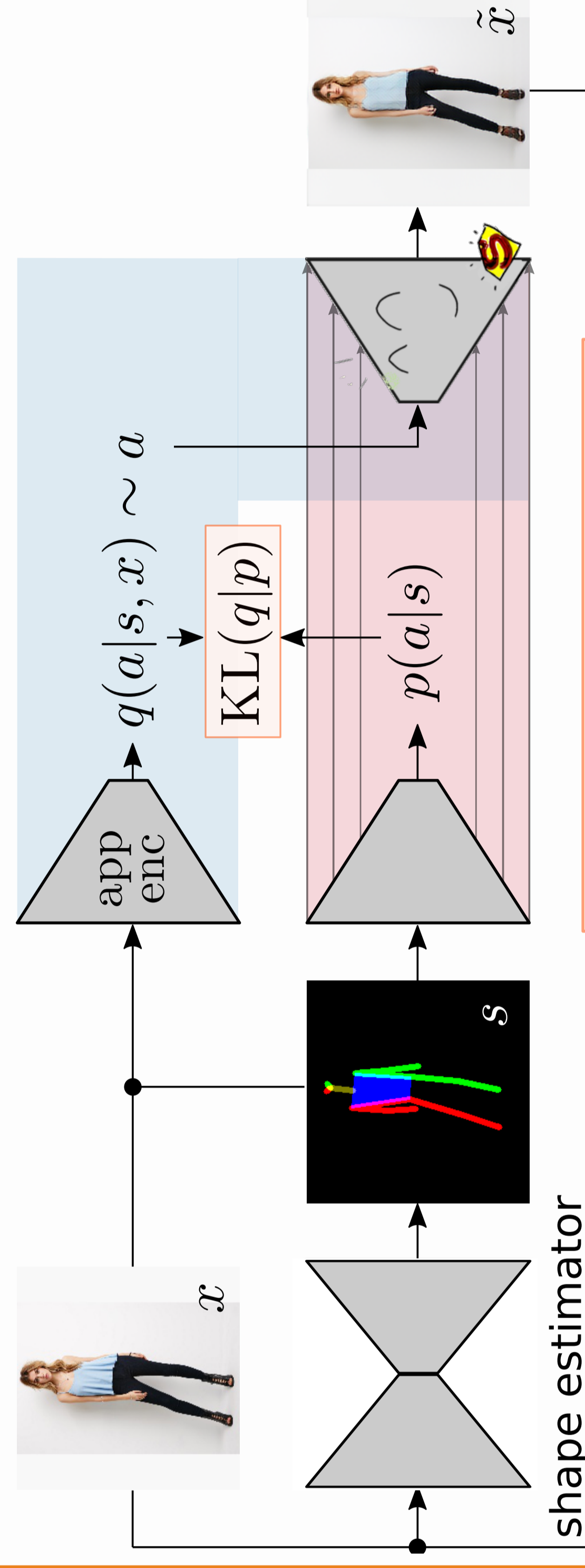
Method

Model images as being generated by two underlying factors: appearance and shape



when modeling both factors as latent variables with a VAE, we cannot disentangle them

\Rightarrow use an estimate of the shape and then learn to disentangle appearance



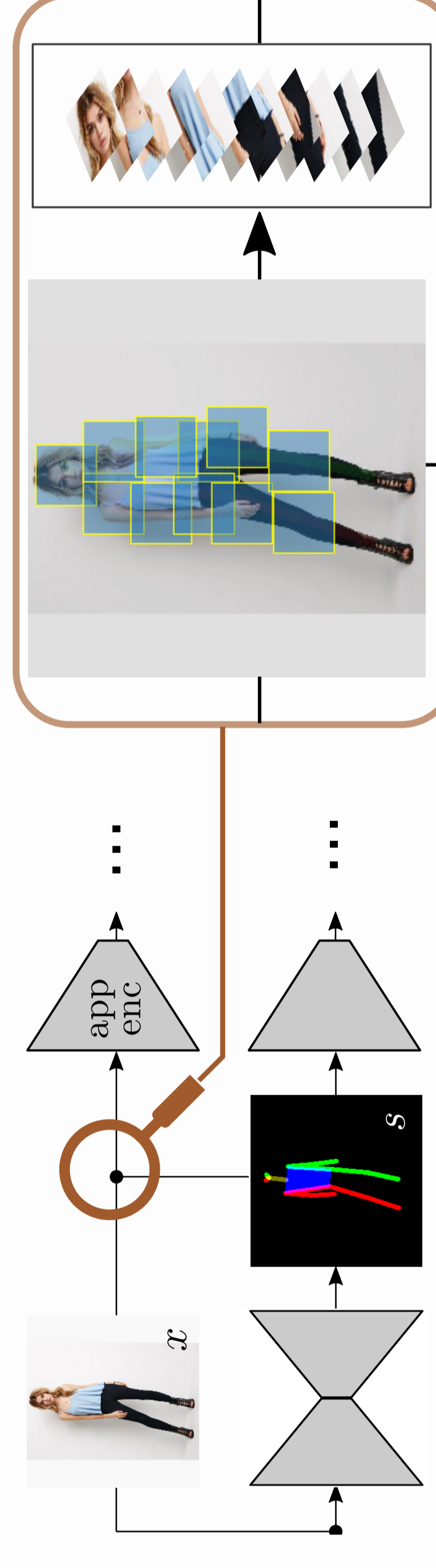
Fit model using shape conditional ELBO:

$$\log p(x|s) \geq \mathbb{E}_{a \sim q(a|x,s)} \log p(x|s,a) - \text{KL}(q(a|x,s) \| p(a|s))$$

$\log p(x|s, a) = -L(x, \tilde{x})$ reconstruction disentanglement

- perceptual loss of VGG19 for reconstruction
- retain spatial information of shape using skip-connections from shape estimate to generator
- KL term penalizes redundancies between appearance and shape \Rightarrow disentangling of appearance

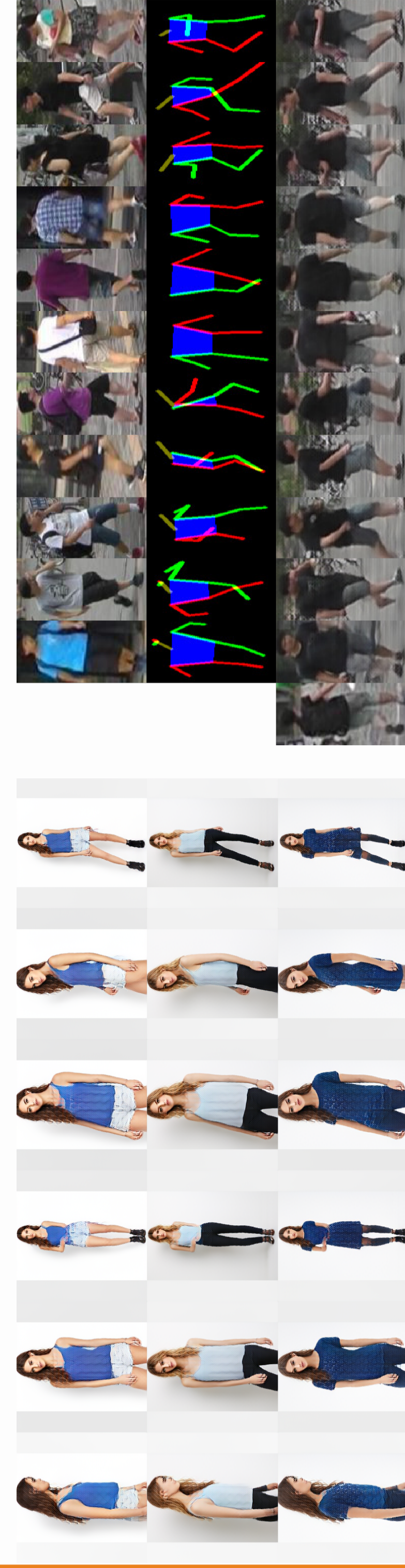
Cropping parts to localize appearance:



crop part boxes and stack them channelwise as input for appearance encoder

Results

Disentangled synthesis



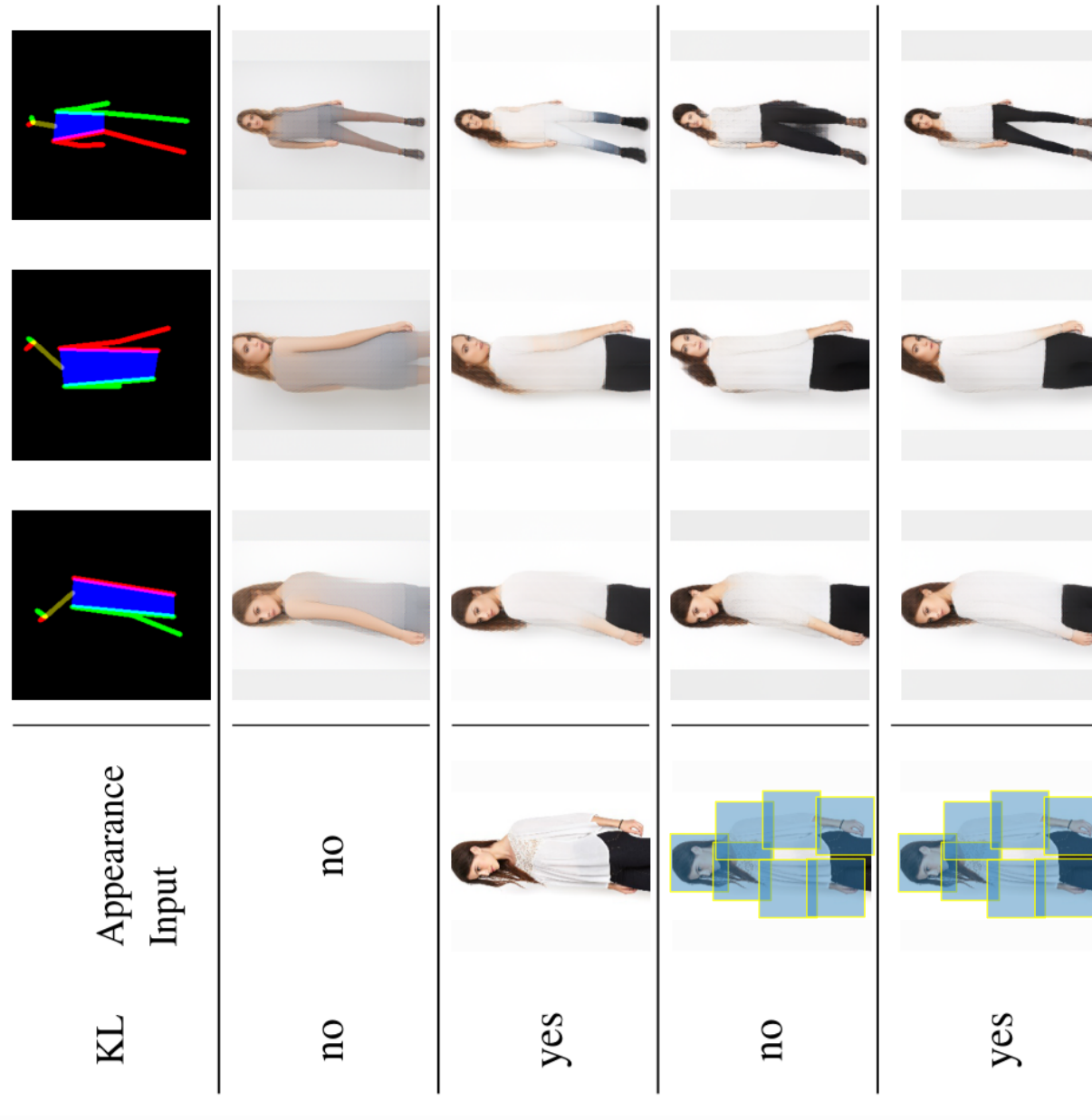
Variations in appearance and shape can be synthesized independently.

Appearance and shape can be transferred between images.

Other sources of shape



Ablation study



Edges can be used as shape estimates, too. Ablation experiments show that disentanglement with both the KL term and object parts is required for highly articulated objects.

Qualitative and quantitative comparisons

method	Market1501			DeepFashion		
	RS	SSIM	IS	RS	SSIM	IS
real data	mean	std	mean	std	mean	std
	3.678	0.274	1.000	0.000	3.415	0.399
PGZ	3.326	0.340	2.668	2.668	0.779	-
PGZ-G+poseMaskedLoss	3.490	0.283	3.094	3.094	0.761	-
PGZ-G+D	3.460	0.253	3.090	3.090	0.762	-
PGZ-G+G2+D	2.289	0.0489	0.166	0.166	2.640	0.2171
our	3.214	0.119	0.353	0.107	3.087	0.2394
					0.786	0.068

dataset	Our		PGZ	
	f6	f6	f6	f6
Market1501	55.95	125.99	67.39	151.16
COCO	23.23	59.26	15.53	69.57
DeepFashion	7.34	133.83	59.24	140.66
Market1501	54.60	59.95	56.24	59.73

Qualitative and quantitative comparisons show improvements over previous methods. (a) SSIM scores are improved, (b) keypoints are more accurately recovered from synthetic images and (c) appearance embeddings are more compact.

Video, more results, code and data
<https://compvis.github.io/vunet/>

