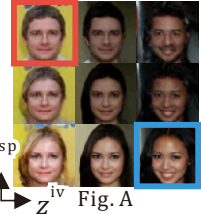


1 We thank the reviewers for fruitful comments. Here, we respond to the major comments. For the minor points like  
 2 presentation issues, we will fix them in the camera-ready (and *we do not mention them here due to space limitation.*)

3 **Reviewer #1** (i) We could not compare our method to [21], because [21] was published in this April and their  
 4 code was not available. (ii) As the reviewer mentioned, a couple of images, say  $c_{\text{smile}}$  and  $c_{\text{non-smile}}$ , were used as  
 5 conditions for smiling and non-smiling images, respectively. Our model is trained using  $\{(x, c_{\text{smile}}) \mid x \in X_{\text{smile}}\}$  and  
 6  $\{(x, c_{\text{non-smile}}) \mid x \in X_{\text{non-smile}}\}$ , where  $X_{\text{smile}}$  and  $X_{\text{non-smile}}$  are sets consisting of smiling and not-smiling images in  
 7 CelebA, respectively. Our model then learns the difference between  $x \in X_{\text{smile}}$  and  $x \in X_{\text{non-smile}}$  in the embedded  
 8 space, in which  $c_{\text{smile}}$  and  $c_{\text{not-smile}}$  are only used to distinguish whether  $x \in X_{\text{smile}}$  or  $x \in X_{\text{non-smile}}$ . In contrast, our  
 9 CHC experiment demonstrates a situation in which the condition image (i.e., the initial microstructure) of each image  
 10 has more diversity. (iii) We used Inception-V3 to calculate FID for CHC. (iv) We apologize that the description “make  
 11 a part of  $z$  follow the non-informative distribution  $\mathcal{N}(0, 1)$ ” in line 125 of our submitted manuscript was incorrect. It  
 12 should be modified as “make a part of  $z$  follow  $\mathcal{N}(f_{\mu}(b_{\mu}), \text{diag}(\exp f_{\sigma}(b_{\sigma})))$  that is independent of  $c$ ”.

13 **Reviewer #2** (i) We trained FUNS using CelebA with Male-Female conditions. An interpolated  
 14 result is shown in Fig. A. We can see that  $z$  is successfully disentangled into the condition  
 15 dependent ( $z^{\text{sp}}$ ) and independent ( $z^{\text{iv}}$ ) parts.



16 (ii) We trained a ResNet to predict smile/non-smile for CelebA, and then, evaluated the prediction  
 17 accuracy for images generated by FUNS, PUNet, and VUNet, as well as the real test set of CelebA.  
 18 For CHC, we trained in-house U-Net to predict the initial condition, and then, evaluated the  $L_2$   
 19 prediction error using  $L_2$  distance between the ground truth and the prediction. The results are  
 20 summarized in Table A. Note that the prediction error for the CHC prediction is 22.84 when the  
 21 ground truth are randomly shuffled. We see that every model can generate images that depend on the conditions in  
 22 terms of the prediction accuracy. PUNet and VUNet achieved higher prediction accuracy compared to Real/FUNS in  
 23 both CelebA and CHC. This may be because PUNet and VUNet generated similar images that are easy to predict the  
 24 condition. As the c-LPIPS scores suggest, PUNet and VUNet generated less diverse images compared to Real/FUNS.

(Table A)	CelebA		CHC	
	Acc.	c-LPIPS	Err.	c-LPIPS
Real	0.924	0.284	9.54	0.169
FUNS	0.973	0.262	9.67	0.157
PUNet	1.000	0.180	9.48	0.108
VUNet	0.977	0.146	9.46	0.118

25 (iii) We have not tried to train FUNS with CelebA-HQ. It is uncertain if FUNS can  
 26 be successfully trained using HQ images (further techniques might be required).  
 27 However, as shown in Table 1 of our paper, FUNS has advantages over VAE-based  
 28 I2I models (incl. PUNet, VUNet) in the image quality (FID) and diversity (LPIPS).  
 29 The main advantage of flow-based models over GAN-based ones is, as discussed  
 30 in [14], that they have invertible mappings between images and latent codes, which  
 31 will be useful for downstream tasks, e.g., Gómez-Bombarelli, R., et al., *ACS Cent. Sci.*, **4**, (2018), 268–276.

32 (iv) We carried out an ablation study to show the effect of each loss term (only with CHC  
 33 dataset due to time limitation). We first trained FUNS using only  $\mathcal{L}^{\text{flow}}$ , and then, other  
 34 loss terms are added sequentially. The results are shown in Table B, in which  $L_2$  errors  
 35 (explained in Table A) and the number of non-zero elements in  $M$  ( $\text{dim}(z^{\text{sp}})$ ) are reported.  
 36 The lower error means that the generated images are more related to the respective conditions.

(Table B)	Err.	$\text{dim}(z^{\text{sp}})$
$\mathcal{L}^{\text{flow}}$	9.56	4,096
$+\mathcal{L}^{\text{recons}}$	9.51	4,040
$+\mathcal{L}^{\text{squeeze}}$	9.58	1,095
$+\mathcal{L}^{\text{entropy}}$	9.67	550

37 From Table B, if we train FUNS by using only  $\mathcal{L}^{\text{flow}}$ ,  $\text{dim}(z^{\text{sp}}) = 4096$ , which is in fact equal to the whole latent  
 38 variable (therefore then  $\text{dim}(z^{\text{iv}}) = 0$ ), suggesting that all the latent variables are dependent on the condition  $c$ . In  
 39 contrast, by adding proposed loss terms,  $\text{dim}(z^{\text{sp}})$  decreases to 550 (then  $\text{dim}(z^{\text{iv}}) = 3546$ ), while the prediction error  
 40 is almost maintained. This result suggests that the input image  $x$  is successfully disentangled in the latent space into  
 41 condition-dependent/independent parts. The impact of dropping  $\mathcal{L}^{\text{entropy}}$  is also illustrated in Fig. 5 of our manuscript.

42 **Reviewer #3** (i) We modified our manuscript to clarify the basic idea in the earlier part  
 43 of the paper according to the reviewer’s comment. (ii) As the reviewer pointed out, it  
 44 is interesting to apply variational inference to train the encoder-decoder in our model,  
 45 which will be described as a future work. (iii) As shown in Table B, 550 of latent  
 46 variables were unmasked (therefore the remaining 3546 latents were finally masked) for  
 47 CHC dataset. In more detail, let  $z^{(l)}$  be the latent variables at level  $l$  ( $z^{(1)} \in \mathbb{R}^{32 \times 32 \times 2}$ ,  
 48  $z^{(2)} \in \mathbb{R}^{16 \times 16 \times 4}$ ,  $z^{(3)} \in \mathbb{R}^{8 \times 8 \times 8}$ , and  $z^{(4)}, z^{(5)} \in \mathbb{R}^{4 \times 4 \times 16}$ ), the numbers of unmasked  
 49 latents were 2, 2, 200, 181, and 165 for  $l = 1, \dots, 5$ , respectively. It means that the  
 50 latents in higher resolution layers are more likely to be masked. (iv) We generated  
 51 CelebA samples using FUNS with the condition of Smile by varying only a part of  $z^{(l)}$ , while the remaining  $z^{(l)}$  are  
 52 fixed to  $f_{\mu}(b_{\mu}^{(l)})$ . Fig. B shows the results: (top)  $z^{(1)}, z^{(2)}$ , and  $z^{(3)}$  are sampled. (middle)  $z^{(4)}$  is sampled. (bottom)  
 53  $z^{(5)}$  is sampled. We can see that larger diversity is captured by latent variables in lower resolution layers ( $z^{(4)}$  and  $z^{(5)}$ ).  
 54 Latents in high resolution layers ( $z^{(1)}, z^{(2)}$ , and  $z^{(3)}$ ) control very subtle facial expressions. (v) Please refer to response  
 55 (ii) to Reviewer 2 for the experimental results that evaluate the proposed model in terms of the prediction accuracy by  
 56 pretrained classifiers. (vi) We will correct grammatical issues through an English proofreading service.

