A Network architecture & training details 1

Human face translation The network architecture used for image translation between *sketch*, 2 photo and paint is given in Table A. 3

Encoder/Decoder utilizes either convolutional layer or transposed convolutional layer, and Discrimi-4

nator utilizes convolutional layer or fully connected layer (FC). We note that the domain vector v is a 5

3 dimensional one-hot vector combined with a 4 dimensional binary vector (2 binary bits for each 6

attribute) for the attribute of interest. (7 dimensions in total) 7

During training, we employ ADAM optimizer with $\beta_1 = 0.5$ and $\beta_2 = 0.99$. The learning rate is 8

set as 1e - 4, and the batch size is set as 96 (32 from each domain). We apply WGAN-GP as our 9

pixel space discriminator. The weight for the objective in Equation 6 is set as follow: 2e - 4 for 10

KL divergence in \mathcal{L}_{vae} , 0.01 for \mathcal{L}_{G}^{adv} and \mathcal{L}_{cls} , and 1 for others. *KL divergence* is ignored in the experiment of Table 2 since generating unseen samples is not necessary. 11

12

Encoder Input : 64x64x3									
Layer	Filters	Kernel size	Stride	BatchNorm.	Activation				
1	64	4x4	2		Leaky ReLU				
1	0.			v					
2	128	4x4	2	V	Leaky ReLU				
3	256	4x4	2	V	Leaky ReLU				
4	512	4x4	2	\checkmark	Leaky ReLU				
5	1024	4x4	2	\checkmark	Leaky ReLU				
μ	1024	4x4	2	-	-				
Generator									
Input : 1024+7									
Layer	Filters	Kernel size	Stride	BatchNorm.	Activation				
1	1024	4x4	2	\checkmark	Leaky ReLU				
2	512	4x4	2	\checkmark	Leaky ReLU				
3	256	4x4	2	\checkmark	Leaky ReLU				
4	128	4x4	2	\checkmark	Leaky ReLU				
5	64	4x4	2	\checkmark	Leaky ReLU				
6	3	4x4	2	-	Tanh				
		Dise	criminat	or					
	Input : 64x64x3								
Layer	Filters	Kernel size	Stride	BatchNorm.	Activation				
1	16	4x4	2	-	Leaky ReLU				
2	32	4x4	2	-	Leaky ReLU				
3	64	4x4	2	-	Leaky ReLU				
4	128	4x4	2	-	Leaky ReLU				
5 (FC)	512	-	-	-	Leaky ReLU				
6 (FC)	[1,7]	-	-	-	-				

Table A: The network architecture of UFDN for *sketch*, *photo* and *paint*.

UDA on digits The network architecture and training hyper-parameters used for $MNIST \rightarrow USPS$ 13 and $USPS \rightarrow MNIST$ are identical as illustrated in Table B. 14

Encoder/Decoder uses convolutional layer/transposed convolutional layer. The digit classifier is a 15

single layer fully-connected network which is jointly learned with our UFDN (only labels in source 16

domain). We note that the domain vector v is a 2 dimensional one-hot vector. we employ ADAM 17

optimizer with $\beta_1 = 0.5$ and $\beta_2 = 0.99$. The learning rate is set as 1e - 4, and the batch size is set as 18

19 32 (16 from each domain).

The weight for the objective in Equation 6 is set as follow: 1e - 7 for *KL divergence* in \mathcal{L}_{vae} , 0.1 for \mathcal{L}_{E}^{adv} and 1 for others. In the task of UDA, we discover that the pixel space discriminator D_x is not necessary since we perform classification on domain-invariant representations instead of images. 20 21

22

The network architecture used for UDA on SVHN \rightarrow MNIST is given in Table C. 23

	Encoder Input : 32x32x3							
Layer	Filters	Kernel size	Stride	BatchNorm.	Activation			
1	64	4x4	2	\checkmark	Leaky ReLU			
2	128	4x4	2	\checkmark	Leaky ReLU			
3	256	4x4	2	\checkmark	Leaky ReLU			
4	512	4x4	2	\checkmark	Leaky ReLU			
μ	1024	4x4	2	-	-			
	Generator							
	Input : 64+2							
Layer	Filters	Kernel size	Stride	BatchNorm.	Activation			
1	512	4x4	2	\checkmark	Leaky ReLU			
2	256	4x4	2	\checkmark	Leaky ReLU			
3	128	4x4	2	\checkmark	Leaky ReLU			
4	64	4x4	2	\checkmark	Leaky ReLU			
5	3	4x4	2	-	Tanh			
Digit Classifier								
Input : 64								
Layer	Filters	Kernel size	Stride	BatchNorm.	Activation			
1 (FC)	10	-	-	-	Soft-max			

Table B: The network architecture of UFDN for $MNIST \rightarrow USPS$ and $USPS \rightarrow MNIST$.

- 24 Encoder/Decoder utilizes either convolutionalal layer or transposed convolutionalal layer, and the
- ²⁵ digit classifier is a single layer fully-connected network that is jointly learned with our UFDN (with
- labels in source domain). We note that the domain vector v is a 2 dimensional one-hot vector.
- During training, we employ ADAM optimizer with $\beta_1 = 0.5$ and $\beta_2 = 0.99$. The learning rate is set as 1e - 4, and the batch size is set as 32 (16 from each domain).

Encoder							
Input : 32x32x3							
Layer	Filters	Kernel size	Stride	BatchNorm.	Activation		
1	128	4x4	2	\checkmark	Leaky ReLU		
2	256	4x4	2	\checkmark	Leaky ReLU		
3	512	4x4	2	\checkmark	Leaky ReLU		
4	1024	4x4	2	\checkmark	Leaky ReLU		
μ	2048	4x4	2	-	-		
Generator							
Input : 2048+2							
Layer	Filters	Kernel size	Stride	BatchNorm.	Activation		
1	1024	4x4	2	\checkmark	Leaky ReLU		
2	512	4x4	2	\checkmark	Leaky ReLU		
3	256	4x4	2	\checkmark	Leaky ReLU		
4	128	4x4	2	\checkmark	Leaky ReLU		
5	3	4x4	2	-	Tanh		
Digit Classifier							
Input : 2048							
Layer	Filters	Kernel size	Stride	BatchNorm.	Activation		
1 (FC)	10	-	-	-	Soft-max		

Table C: The network architecture of UFDN for SVHN \rightarrow MNIST.