

Iterative Feature Transformation for Fast and Versatile Universal Style Transfer

Anonymous ECCV submission

Paper ID 3170

This document supplements Sections 3 and 4 of the main paper. In particular, it includes the following:

- Derivation of the analytical gradient (supplements **Section 3.2**).
- Training details of the autoencoders (supplements **Section 4**).
- Stylized results for quantitative analysis of photo-realistic transfer (supplements **Section 4.2**).
- Formulation of NST and WCT for multi-style transfer and double-style transfer results from AdaIN and Avatar-net (supplements **Section 4.3**).

1 Derivation of the analytical gradient

For simplicity, we suppress the subscript N . Here we show that if

$$l_j(\mathbf{F}) = \|\mathbf{F} - \mathbf{F}^{(j)}\|_F^2 + \lambda \left\| \frac{1}{n} \mathbf{FF}^T - \frac{1}{m} \mathbf{F}_s \mathbf{F}_s^T \right\|_F^2, \quad (1)$$

then

$$\frac{dl}{d\mathbf{F}} = 2(\mathbf{F} - \mathbf{F}^{(j)}) + \frac{4\lambda}{n} \left(\frac{1}{n} \mathbf{FF}^T - \frac{1}{m} \mathbf{F}_s \mathbf{F}_s^T \right) \mathbf{F}. \quad (2)$$

Proof.

$$\|\mathbf{F} - \mathbf{F}^{(j)}\|_F^2 \quad (3)$$

$$= \text{tr}[(\mathbf{F} - \mathbf{F}^{(j)})(\mathbf{F} - \mathbf{F}^{(j)})^T] \quad (4)$$

$$= \text{tr}[\mathbf{FF}^T - 2\mathbf{F}(\mathbf{F}^{(j)})^T + \mathbf{F}^{(j)}(\mathbf{F}^{(j)})^T], \quad (5)$$

and

$$\left\| \frac{1}{n} \mathbf{FF}^T - \frac{1}{m} \mathbf{F}_s \mathbf{F}_s^T \right\|_F^2 \quad (6)$$

$$= \text{tr}\left[\left(\frac{1}{n} \mathbf{FF}^T - \frac{1}{m} \mathbf{F}_s \mathbf{F}_s^T\right)\left(\frac{1}{n} \mathbf{FF}^T - \frac{1}{m} \mathbf{F}_s \mathbf{F}_s^T\right)\right] \quad (7)$$

$$= \text{tr}\left[\frac{1}{n^2} \mathbf{FF}^T \mathbf{FF}^T - \frac{2}{nm} \mathbf{FF}^T \mathbf{F}_s \mathbf{F}_s^T + \frac{1}{m^2} \mathbf{F}_s \mathbf{F}_s^T \mathbf{F}_s \mathbf{F}_s^T\right]. \quad (8)$$

Let $\mathbf{F} = [f_1, f_2, \dots, f_n]$, $\mathbf{F}^{(j)} = [f_1^{(j)}, f_2^{(j)}, \dots, f_n^{(j)}]$, and $\mathbf{F}_s = [f_1^s, f_2^s, \dots, f_m^s]$. We first find the partial derivatives with respect to f_i :

$$\frac{\partial \|\mathbf{F} - \mathbf{F}^{(j)}\|_F^2}{\partial f_i} = \frac{\partial \text{tr}[\mathbf{FF}^T]}{\partial f_i} - 2 \frac{\partial \text{tr}[\mathbf{F}(\mathbf{F}^{(j)})^T]}{\partial f_i}, \quad (9)$$

$$\frac{\partial \|\frac{1}{n}\mathbf{F}\mathbf{F}^T - \frac{1}{m}\mathbf{F}_s\mathbf{F}_s^T\|_F^2}{\partial f_i} = \frac{1}{n^2} \frac{\partial \text{tr}[\mathbf{F}\mathbf{F}^T\mathbf{F}\mathbf{F}^T]}{\partial f_i} - \frac{2}{nm} \frac{\partial \text{tr}[\mathbf{F}\mathbf{F}^T\mathbf{F}_s\mathbf{F}_s^T]}{\partial f_i}, \quad (10)$$

where $\text{tr}[\mathbf{F}\mathbf{F}^T] = \sum_{a=1}^n f_a^T f_a$, $\text{tr}[\mathbf{F}(\mathbf{F}^j)^T] = \sum_{a=1}^n f_a^T f_a^{(j)}$,

$$\text{tr}[\mathbf{F}\mathbf{F}^T\mathbf{F}\mathbf{F}^T] \quad (11)$$

$$= \text{tr}[\sum_{a=1}^n f_a f_a^T \sum_{b=1}^n f_b f_b^T] \quad (12)$$

$$= \sum_{a=1}^n \sum_{b=1}^n \text{tr}[f_a f_a^T f_b f_b^T] \quad (13)$$

$$= \sum_{a=1}^n \sum_{b=1}^n (f_a^T f_b)^2, \quad (14)$$

and similar to $\text{tr}[\mathbf{F}\mathbf{F}^T\mathbf{F}\mathbf{F}^T]$, we have $\text{tr}[\mathbf{F}\mathbf{F}^T\mathbf{F}_s\mathbf{F}_s^T] = \sum_{a=1}^n \sum_{b=1}^m (f_a^T f_b^s)^2$.

For the partial derivatives with respect to f_i , we only have to focus on the terms associated with f_i . Therefore,

$$\frac{\partial \text{tr}[\mathbf{F}\mathbf{F}^T]}{\partial f_i} = \frac{\partial f_i^T f_i}{\partial f_i} = 2f_i, \quad (15)$$

$$\frac{\partial \text{tr}[\mathbf{F}(\mathbf{F}^{(j)})^T]}{\partial f_i} = \frac{\partial f_i^T f_i^{(j)}}{\partial f_i} = f_i^{(j)}, \quad (16)$$

$$\frac{\partial \text{tr}[\mathbf{F}\mathbf{F}^T\mathbf{F}\mathbf{F}^T]}{\partial f_i} \quad (17)$$

$$= \frac{\partial}{\partial f_i} \left(\sum_{a \neq i} (f_a^T f_i)^2 + \sum_{b \neq i} (f_i^T f_b)^2 + (f_i^T f_i)^2 \right) \quad (18)$$

$$= 2 \sum_{a \neq i} (f_a^T f_i) f_a + 2 \sum_{b \neq i} (f_i^T f_b) f_b + 4(f_i^T f_i) f_i \quad (19)$$

$$= 4 \sum_{a \neq i} (f_a^T f_i) f_a + 4(f_i^T f_i) f_i \quad (20)$$

$$= 4 \sum_{a=1}^n (f_a^T f_i) f_a \quad (21)$$

$$= 4\mathbf{F}\mathbf{F}^T f_i, \quad (22)$$

and

$$\frac{\partial \text{tr}[\mathbf{F}\mathbf{F}^T\mathbf{F}_s\mathbf{F}_s^T]}{\partial f_i} = \frac{\partial}{\partial f_i} \sum_{b=1}^m (f_i^T f_b^s)^2 = 2 \sum_{b=1}^m (f_i^T f_b^s) f_b^s = 2 \sum_{b=1}^n ((f_b^s)^T f_i) f_b^s = 2\mathbf{F}_s\mathbf{F}_s^T f_i. \quad (23)$$

Putting everything together, we have

$$\frac{\partial l_j(\mathbf{F})}{\partial f_i} = 2f_i - 2f_i^{(j)} + \lambda(4\frac{1}{n^2}\mathbf{FF}^T f_i - 4\frac{1}{nm}\mathbf{F}_s\mathbf{F}_s^T f_i) \quad (24)$$

$$= 2(f_i - f_i^{(j)}) + \frac{4\lambda}{n}(\frac{1}{n}\mathbf{FF}^T - \frac{1}{m}\mathbf{F}_s\mathbf{F}_s^T)f_i. \quad (25)$$

Finally,

$$\frac{dl_j(\mathbf{F})}{d\mathbf{F}} = \left[\frac{\partial l_j}{\partial f_1}, \frac{\partial l_j}{\partial f_2}, \dots, \frac{\partial l_j}{\partial f_n} \right] \quad (26)$$

$$= 2(\mathbf{F} - \mathbf{F}^{(j)}) + \frac{4\lambda}{n}(\frac{1}{n}\mathbf{FF}^T - \frac{1}{m}\mathbf{F}_s\mathbf{F}_s^T)\mathbf{F}. \quad (27)$$

2 Training details of the autoencoders

The four autoencoders are trained by minimizing an image reconstruction loss and a perceptual loss. In particular, if the functions of the $encoder_N$ and $decoder_N$ are denoted $\phi_N(\cdot)$ and $\psi_N(\cdot)$, respectively, the $decoder_N$ is trained by minimizing the loss \mathcal{L}_{AE} :

$$\mathcal{L}_{AE} = ||\psi_N(\phi_N(I)) - I||_F^2 + ||\phi_N(\psi_N(\phi_N(I))) - \phi_N(I)||_F^2, \quad (28)$$

where I is an input image. We train the autoencoders on the MS-COCO dataset. To support batch training, each image from the dataset is resized to 512×512 and randomly cropped to 256×256 as a training example in a batch. For the autoencoders associated with $relu4_1$ and $relu3_1$ layers, they are trained with a batch size of 8 for 5 epochs, while for $relu2_1$ and $relu1_1$ cases, the autoencoders are trained for 3 epochs, due to their smaller sizes. We use Adam optimizer with the learning rate 1×10^{-4} and without weight decay. Moreover, we use up-sampling layers with bilinear interpolation in the decoders as the symmetric part of the max-pooling layers in the encoders.

3 Stylized results for quantitative analysis of photo-realistic transfer

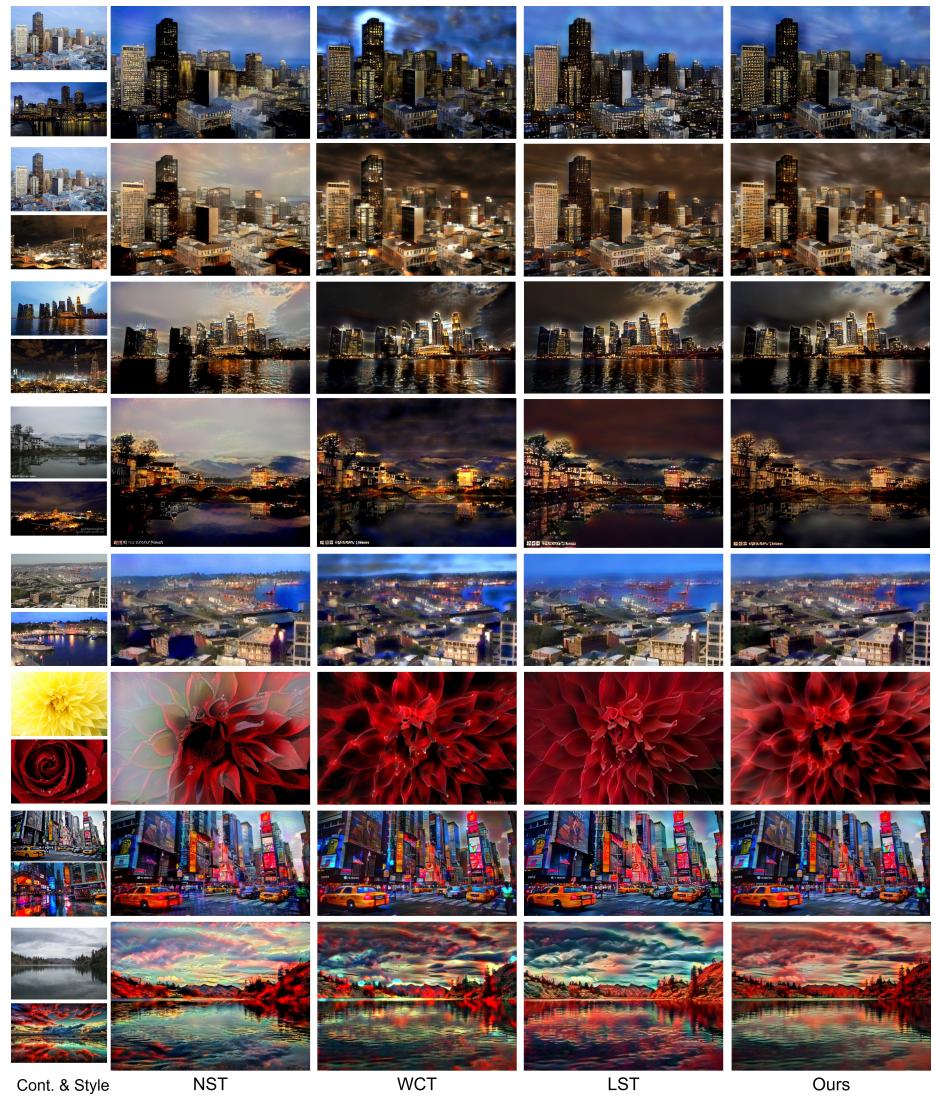


Fig. 1: Photo-realistically stylized images from 30 pairs of a content and a style images for quantitative analysis. No spatial control and no post-processing are applied (Part 1/4).

180
181
182
183

Table 1: Speed performance of our method under $n_{upd} = 15$ and $n_{iter} = 1$ for generating the results in figures 1, 2, 3, and 4. **Unit:** Second.

	256 × 256	512 × 512	768 × 768	1024 × 1024
time	0.13	0.31	0.62	0.92

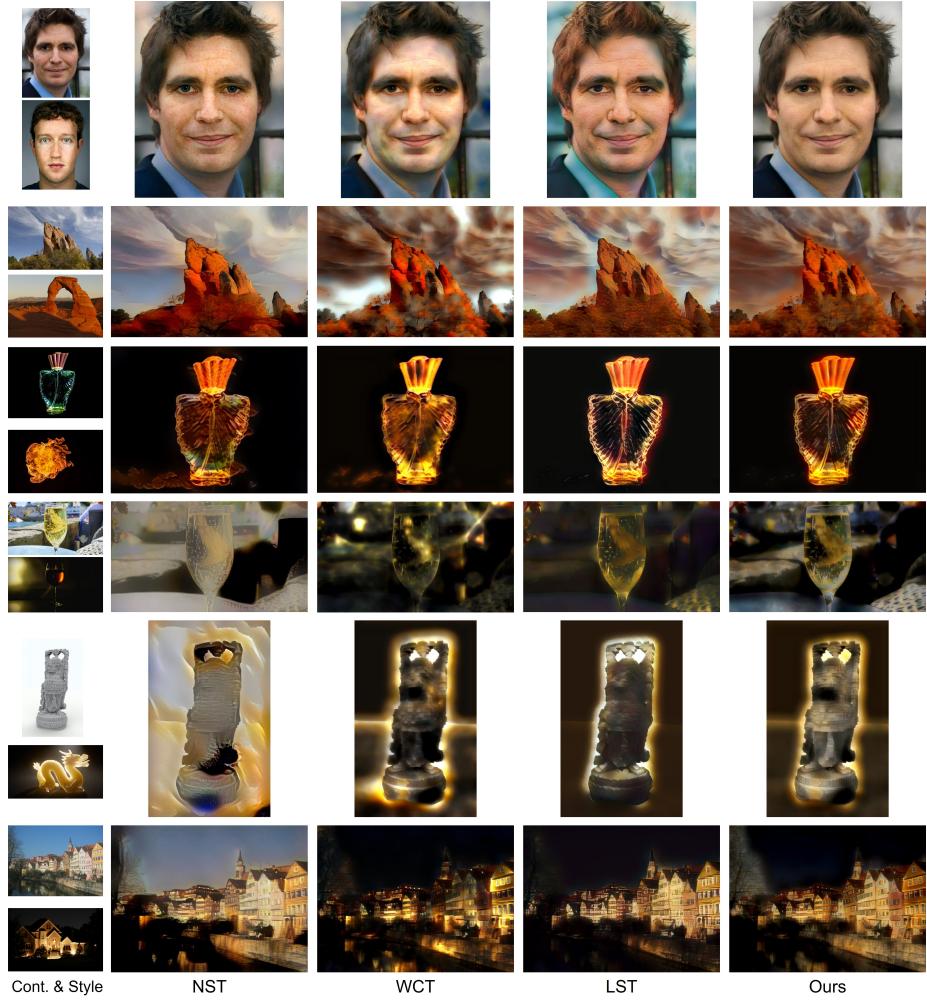
187
188

Fig. 2: Photo-realistically stylized images from 30 pairs of a content and a style images for quantitative analysis. No spatial control and no post-processing are applied (Part 2/4).

224

180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224

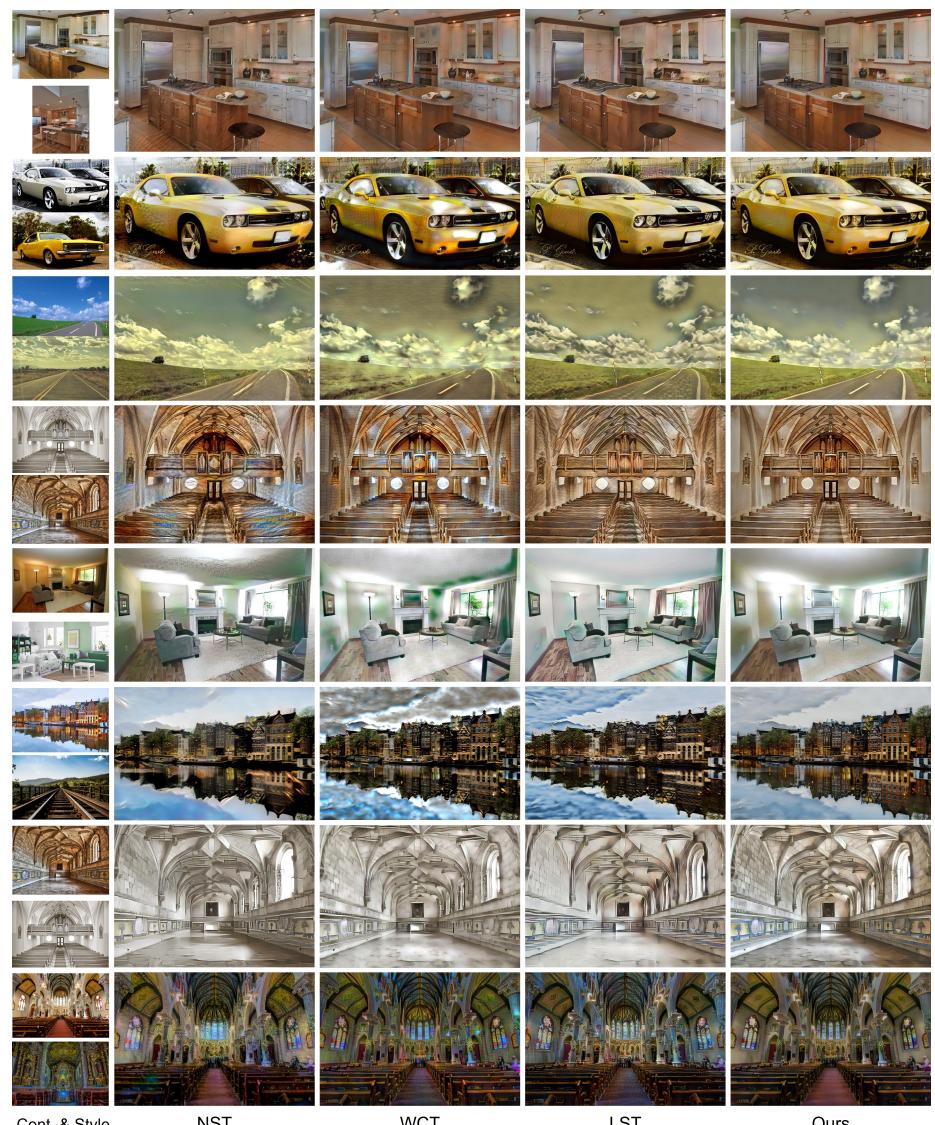


Fig. 3: Photo-realistically stylized images from 30 pairs of a content and a style images for quantitative analysis. No spatial control and no post-processing are applied (Part 3/4).

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314

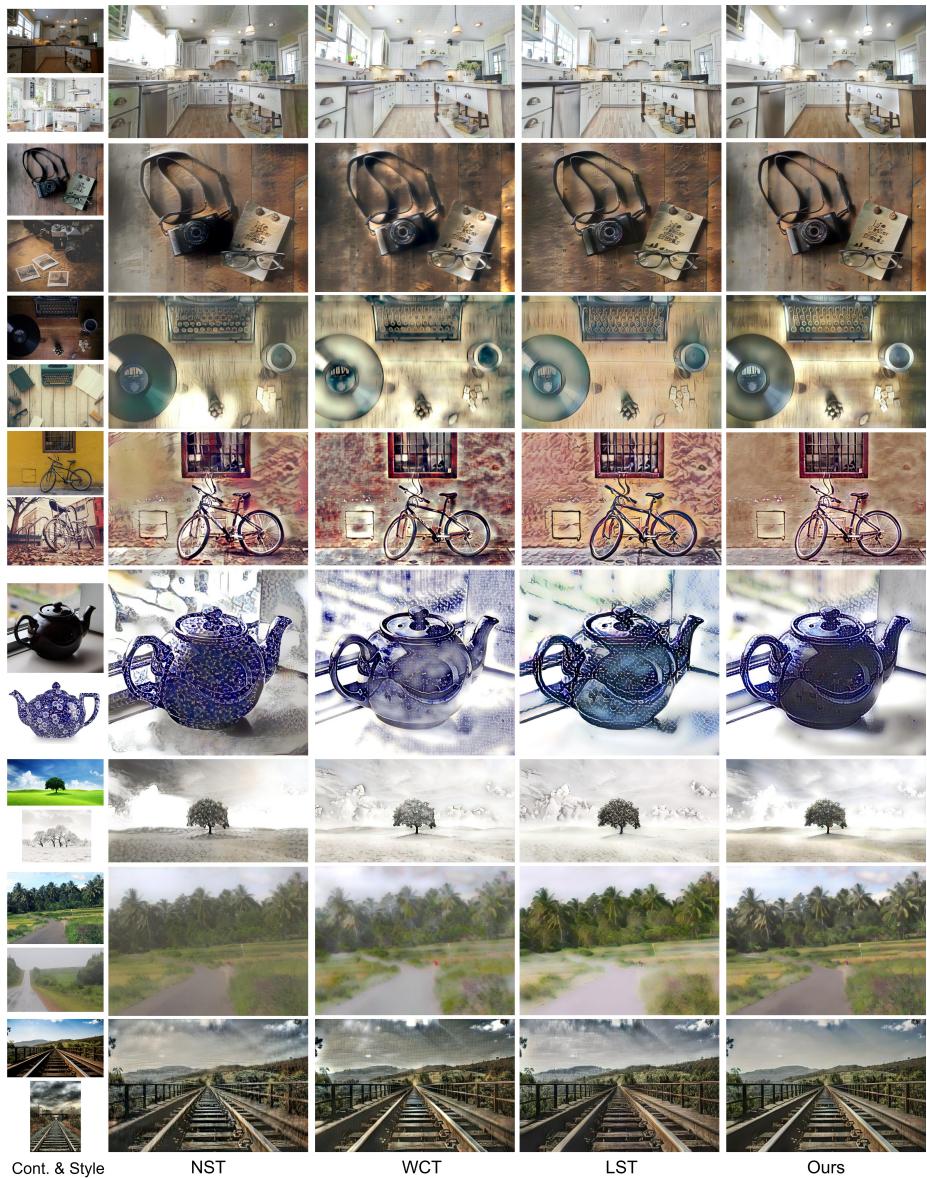


Fig. 4: Photo-realistically stylized images from 30 pairs of a content and a style images for quantitative analysis. No spatial control and no post-processing are applied (Part 4/4).

315 4 Formulation of NST and WCT for multi-style transfer

316 The objective of NST for multi-style transfer is as follows:

$$317 \min_I \|\mathbf{F}_4(I) - \mathbf{F}_{4,c}\|_F^2 + \sum_{N=1}^4 \sum_{k=1}^q \lambda_N^k \left\| \frac{1}{n_N} \mathbf{F}_N(I) \mathbf{F}_N(I)^T - \frac{1}{m_N^k} \mathbf{F}_{N,s}^k (\mathbf{F}_{N,s}^k)^T \right\|_F^2, \\ 318 \quad (29)$$

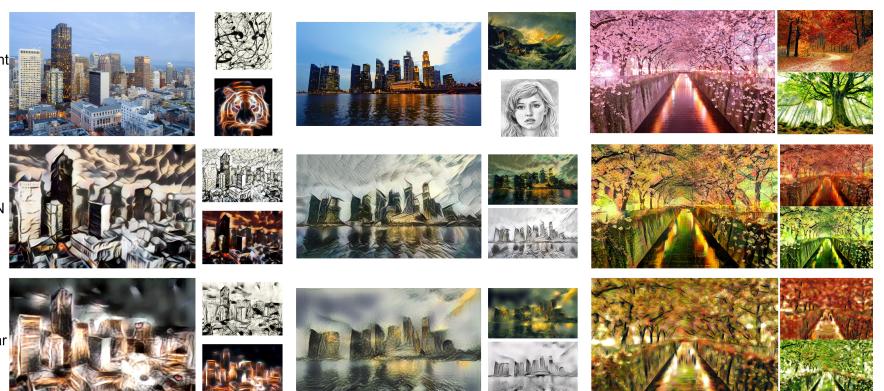
319 where $\mathbf{F}_{N,s}^k$'s are the feature maps of q style images extracted from $encoder_N$.
320 The stylized image is then derived by solving equation 29 using gradient descent
321 by back-propagation. How different style features are included in equation 29 is
322 non-linear.

323 On the other hand, WCT realizes multiple-style transfer by linear interpolation
324 of transformed features. By applying WCT to each style feature $\mathbf{F}_{N,s}^k$ and
325 the content feature $\mathbf{F}_{N,c}$, we can derive a transformed feature $\mathbf{F}_{N,wct}^k$. The final
326 feature $\mathbf{F}_{N,wct}$ to be decoded is an affine combination:

$$327 \mathbf{F}_{N,wct} = \sum_{k=1}^q w_k \mathbf{F}_{N,wct}^k, \text{ with } \sum_{k=1}^q w_k = 1. \quad (30)$$

328 As such, each style is weakened due to $w_k < 1$ in the stylized image and could
329 even not be observed.

330 5 Double-style transfer results from AdaIN and 331 Avatar-net



332 **Fig. 5:** Double-style transfer results from AdaIN and Avatar-net. Unlike our method
333 that preserves the integrity of each style, styles in doubly stylized images from AdaIN
334 and Avatar-net might be weakened due to the linear interpolation of feature maps.