

<http://fit.vut.cz>

## Image Style Transfer

Šimon Sedláček, Marek Sarvaš, Martin Osvald



### 1. Introduction

Image style transfer is simply the problem of applying the stylistic elements of one *style* image to a different, the so-called *content* image. Ideally, the resulting image would look as though it were painted/created by the same artist or technique of the style image, retaining its color scheme, stroke and detail granularity, etc. In this project, we explore and combine two novel approaches to style transfer, ArtFlow [1] and AdaAttN [8]. The code for this project was based on the official implementation of ArtFlow and the repository can be found on <https://github.com/MarekSarvas/ArtFlow>.

### 2. Neural Style Transfer

The first method that has used neural networks to achieve painting style transfer was presented in [3]. The authors of this paper used a pre-trained VGG convolutional neural network to separate the style and content factors from the input image. The content factor can be thought of as 'everything relevant for detecting and classifying objects present in the image'. On the other hand, the style factor could be regarded as a noise factor, that is spatially global to the whole image and has no effects on the classification/detection results. The exact principle of extracting the the content and style factors of an image using the VGG network can be found in [3], we will suffice with denoting this process as follows:

$$F_c = E_c(I) \quad (1)$$

$$F_s = E_s(I), \quad (2)$$

where  $F_c$  and  $F_s$  denote the extracted content and style feature maps, respectively,  $E_c$  and  $E_s$  denote the content and style extraction passes through the decoder, respectively and  $I$  is the input image.

Once separated, the content and style feature maps are then used as reference points for the style transfer itself, which in [3] was done by gradually (using gradient descent optimization) synthesising a completely new image from white noise, that would in the end match both the content and the style feature maps, obtained from the input images.

This method is, however, quite time consuming and therefore, other approaches to the style transfer process itself quickly followed. These newer approaches often use the same method for content/style separation. The new addition is, however, a more efficient style transfer module, that simply combines the style and content feature maps in a way that, once again, creates a new image with the style and content of the input images. These approaches include most notably the Adaptive Instance Normalization (AdaIN) approach [4] and the Whitening and Coloring Transform (WCT) [6]. However, there are some newer methods trying to expand on the basic idea of these style transfer modules such as AvatarNet or AdaAttN [10, 8], bringing some interesting results.

### 3. Artflow

Even though most neural style transfer approaches use a CNN trained for object detection to obtain the desired style and content feature maps, there are some novel approaches, which experiment with replacing this part

46 of the style transfer system. One of such approaches  
47 called ArtFlow was recently presented in [1].

48 This method replaces the original VGG-based en-  
49 coder with a new neural network architecture, derived  
50 from the Glow [5] generative normalizing flow. The  
51 main idea behind using such an architecture is its ca-  
52 pability for lossless decoding. This means, that once  
53 we obtain a latent vector  $z$  for an input image, we  
54 can pass this vector backwards through the network to  
55 reconstruct the original image perfectly.

56 This network architecture was used with hopes of  
57 mitigating some of the problems that come with the  
58 default 'inverse VGG' decoder used in other meth-  
59 ods, mainly the problem of the so-called *content leak*.  
60 This flow-based network architecture should help re-  
61 tain more of the content information from the original  
62 image, therefore the stylized result should be less dis-  
63 torted, the content better distinguishable, etc.

### 64 3.1 ArtFlow network architecture

65 The basic building block of the network used in [1] is  
66 the so-called *flow step*. A flow step consists of three  
67 layers, that perform fully reversible transformations  
68 of their input. These layers include a reversible  $1 \times 1$   
69 convolution, and ActNorm layer and an affine coupling  
70 layer. Together, these three layers form one flow step.

71 These flow steps can be then grouped up, creating  
72 one *flow block*. Each flow block is also preceded with a  
73 special *squeeze* operation, which is simply a reversible  
74 alternative to the regular spacial pooling layers of typ-  
75 ical CNNs. The inner workings and also problems  
76 linked with this layer are further discussed in section  
77 5.

78 In contrast to the VGG-based approaches, the Art-  
79 Flow network used in the same manner for both style  
80 and content image encoding – both input images will  
81 go through the entire network, giving us their represen-  
82 tations in the latent space at the output of the network.  
83 These latent  $z$  vectors then serve as input of the used  
84 style transfer module (e.g. AdaIN), resulting in a styl-  
85 ized  $z_{cs}$  latent vector that can then be passed through  
86 the network in reverse to obtain the final stylized im-  
87 age.

88 In the ArtFlow paper, best results were achieved  
89 when employing two blocks of eight flow modules  
90 each in the final model. The whole ArtFlow system is  
91 depicted in figure 1.

### 92 3.2 ArtFlow style transfer

As it was stated in the previous section, we first obtain  
the  $z$ -space representations for both the content and

style images  $I_c$  and  $I_s$ :

$$z_c = A(I_c) \quad (3)$$

$$z_s = A(I_s), \quad (4)$$

with  $A$  denoting the ArtFlow encoder.

These vectors are then passed to the style transfer  
module, e.g. AdaIN. The style transfer module applies  
a rather simple transform to these vectors, resulting in  
a stylized feature map  $z_{cs}$ :

$$z_{cs} = AdaIN(z_c, z_s). \quad (5)$$

This stylized  $z$ -vector is passed backwards through  
the encoder, resulting in the final stylized image:

$$I_{cs} = A^{-1}(z_{cs}). \quad (6)$$

The whole projection-transfer-reversion scheme is  
depicted in figure 1.

### 3.3 Training ArtFlow

It is now obvious that the aim when training the Art-  
Flow system is to find the parameters of the encoder so  
that the style transfer module (which has no learnable  
parameters) can do the best job possible.

In order to do that, some cost functions have to be  
defined for this purpose. There are two const functions  
used, one supervising the content and one supervising  
the style factor of the resulting images:

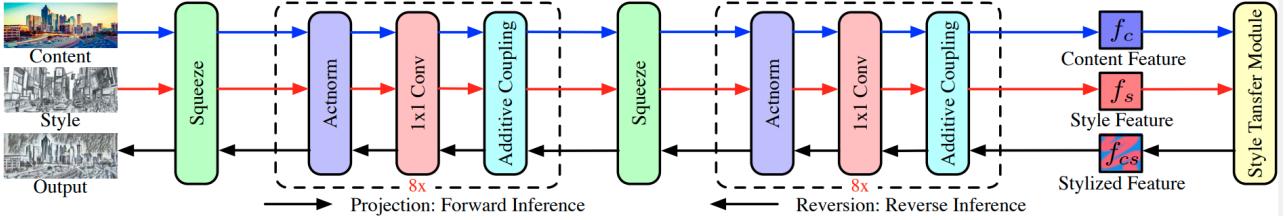
$$L_c = \|A(A^{-1}z_{cs}, z_{cs})\|_2 \quad (7)$$

$$L_s = \sum_{i=1 \dots L} \|\mu(E_{si}(A^{-1}(z_{cs}))) - \mu(E_{si}(I_s))\|_2 \quad (8)$$
$$+ \sum_{i=1 \dots L} \|\sigma(E_{si}(A^{-1}(z_{cs}))) - \sigma(E_{si}(I_s))\|_2.$$

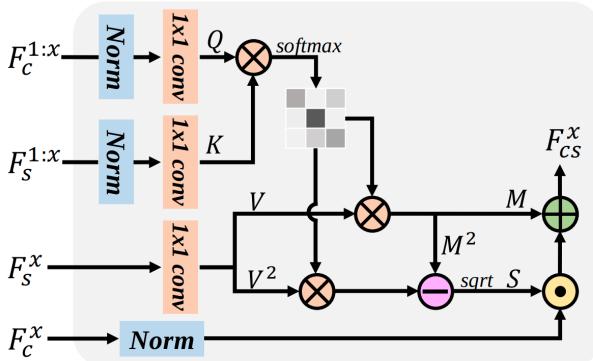
It is important to note that the style loss function  
 $L_s$  makes use of the VGG encoder as a reference point  
for style separation, where  $E_{si}$  denotes a layer in the  
VGG net from which we extract the style features. In  
other words, the whole VGG encoder is a part of the  
loss function.

### 3.4 Style transfer bias

Last important thing to note about the ArtFlow ap-  
proach is that it uses style transfer methods that are *un-  
biased*. This means that the style transfer module will  
perform style transfer without prioritizing content or  
style in the resulting image, retaining maximum origi-  
nal content information. This can be demonstrated by  
stylizing an input content image with a particular style  
and the stylizing this new image again with the style  
of the original input image. ArtFlow achieves perfect



**Figure 1.** Depiction of the whole ArtFlow system with two Glow blocks of eight flow steps each. The diagram also shows the inference directions for both encoding and decoding. Diagram taken from [1].



**Figure 2.** A diagram of the AdaAttN style transfer module taken from [8]. The *Norm* blocks are regular channel-wise mean-variance normalization layers.

reconstruction of the original image this way. More on this topic can be found in section 7.3.

However, we found that some of the results may look 'too clean' and surgical and that some of the biased style transfer modules such as Avatar-Net or AdaAttN [10, 8] achieve subjectively better style transfer by more or less biasing towards style.



**Figure 3.** A showcase of the grid-shaped artifacts present in the stylized images.

block and flow steps within.

146

## 5. Challenges of the proposed method

Even though the qualitative results of our proposed system are satisfactory and comparable to the original ArtFlow results, there are some significant drawbacks that hinder the usefulness of our method in the real world. We will discuss these drawbacks here, before presenting the results themselves, as it is important to understand some of the terms and aspects of the final models that were implemented in reaction to these problems.

147

148  
149  
150  
151  
152  
153  
154  
155  
156

### 5.1 Grid artifacts

The first and probably the biggest limitation of the ArtFlow-based methods are the inevitable grid-shaped artifacts that will appear in the reconstructed images. This phenomenon was not discussed in [1] the original paper at all and posed an unpleasant surprise for us when we got our first results. An example of the grid artifacts on a zoomed-in image can be seen in figure ??

157

158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173

This phenomenon was, however, mentioned in [9] and is caused by the usage of the previously described squeeze layers at the beginning of each flow block. This squeeze layer performs a spatial splitting operation, which essentially splits the original image into multiple channels simply by dissecting neighbouring pixels from each other in a checkerboard-like pattern. Each squeeze layer will therefore quadruple the cur-

126

## 4. Our approach

127

The fact that biased style transfer modules can lead to better stylistic results prompted us to create our own system, which would combine the properties of the ArtFlow encoder and the capabilities of novel style transfer modules such as AdaAttN.

128

AdaAttN is essentially a generalization of AdaIN, taking the same base principle but adding learnable parameters and exploiting attention to better retain the aspects of the original style image. A diagram of the module can be seen in figure 2.

129

In combination with ArtFlow, we use the obtained  $z_c$  vector in place of the  $F_c^x$  and  $F_c^{1,x}$  AdaAttN inputs and the  $z_s$  vector in place of  $F_s^x$  and  $F_s^{1,x}$ .

130

Ultimately, we would like the new system to better capture spacial style aspects than the AdaIN-based baseline, resulting in more lively, genuine and aesthetic stylization.

131

We adopt the same loss functions as regular ArtFlow, experimenting with different amounts of flow

132

174 rent number channels and will half the current width  
175 and height of the image. However, it appears that  
176 after applying style transfer, some some of the local  
177 correlations between pixels are corrupted, resulting in  
178 this grid-like artifact pattern in the final reconstructed  
179 image.

180 The same paper also suggests that a way to mitigate  
181 this phenomenon is to add a few *transition* layers to the model.  
182 Each transition layer is simply the original flow step, but with the affine coupling layer  
183 removed. Adding these transition layers to the model  
184 right after the squeeze operation should allow the  
185 model to learn a linear interpolation between the neighbouring pixels, allowing the model to fight the artifacts.  
186

187 Unfortunately, we had no success with removing  
188 the grid patterns from the images even after adding  
189 the transition layers. It is, however, possible, that  
190 by designing a custom loss function, which would  
191 penalize the grid artifacts in the final image, we would  
192 be able to push the final model into the right direction.  
193

194 What does help is using fewer squeeze operations  
195 in the model. The models with only one flow block  
196 instead of two have the grid effect noticeably reduced,  
197 this is, however, at the cost of slight color saturation  
198 loss and slight general style transfer capability deteriora-  
199 tion. Also, the artifacts are much more visible in im-  
200 ages with many colors. When going for a drawing-type  
201 monochrome stylization, the artifacts tend to recede.  
202

203 The grid pattern is most noticeable on the edges  
204 of each image, resulting not only in strong and visible  
205 grid lines, but also in some unexpected color stripes  
206 along the edges of the image. This is depicted in figure  
207 4. To work around this, we tried to wrap the input  
208 content image with 64 pixels of padding on each side,  
209 effectively trying to push these edge artifacts beyond  
210 the edges of the image before cropping it back to the  
211 original size. This yielded some satisfying results, sig-  
212 nificantly reducing the corruption of the image edges  
213 after style transfer.

## 213 5.2 Memory consumption

214 Because the AdaAttN style transfer module uses the  
215 attention mechanism, it is at one point necessary to  
216 generate an attention map between the key and the  
217 query. This is done by applying the outer product to  
218 these two tensors. However, if these two input query  
219 and key tensors are too large, the outer product matrix  
220 will be quite large.

221 This problem gets even more pronounced when  
222 working with higher resolution images, say 512 pixels  
223 in each dimension. The query and key tensors can then  
224 have shapes that are simply too large for us to compute  
225 their outer product (in some instances, this results in

tensors of more than 40 GB of CUDA memory in size).  
226

227 This issue is even more pronounced when using an  
228 ArtFlow architecture with only one block – the output  
229 vector  $z$  will have most of its original information re-  
230 tained in the original width and height dimensions. As  
231 we increase the number of flow blocks of the network,  
232 more and more of the original image dimensionality  
233 will be moved to the channel dimensions, resulting in  
234 a more manageable width and height values.  
235

236 There is, however, one way to tackle this problem,  
237 which is simply to limit the width and height values of  
238 the query and key tensors, for example by randomly  
239 selecting a number of indices in the original latent  
240 vector to create a new query/key tensor. We observed,  
241 that it is useful to keep the size of this newly created  
242 tensor to more than 64 values per dimension, otherwise  
243 the attention map simply will not retain enough detail  
244 and information to perform well.  
245

246 Even though that limiting the attention map size  
247 does solve the problem partially, most style transfer  
248 instances will still require gigabytes of CUDA memory  
249 to generate the attention maps, making this method  
250 hardly usable in most real life scenarios, at least in its  
251 current state and without a GPU.  
252

## 250 6. Experimental setup

251 For training, we used the MS-COCO [7] dataset as a  
252 source of content images, totalling at 120 thousand  
253 images. For styles, we used the Wikiart<sup>1</sup>, totalling at  
254 40 thousand style images.  
255

256 All models were trained for a total of 40000 iter-  
257 ations with a batch size of 4. The initial learning rate  
258 was set to  $1 \cdot 10^{-4}$  with a decay factor of  $5 \cdot 10^{-5}$ .  
259

260 We make use of the computational resources pro-  
261 vided by Metacentrum VO<sup>2</sup>, specifically the Nvidia  
262 A40 48 GB GPUs on the new Galdor cluster, each  
263 model taking seven to eight hours to train on average.  
264

## 262 7. Results, architecture comparison

263 We train the system in multiple configurations, namely  
264 1x16, 2x8, 1x24, 3x4 flow steps with all of these con-  
265 figurations also having its variant with and without  
266 one or two transition flow layers in each block. In  
267 this section, we will discuss the results of the best per-  
268 forming trained models, as some of the models were  
269 completely unusable – basically any architecture with  
270 three or more flow blocks. The images that such sys-  
271 tems would produce failed to retain enough content  
272 information to be aesthetically pleasing and would

<sup>1</sup><https://www.wikiart.org/>

<sup>2</sup><https://metavo.metacentrum.cz/>



**Figure 4.** Comparison of style transfer on content image (left) and content image padded with 64 pixels with value 0 on each side (right). Padding solves quite noticeable artifacts that appear mostly on the sides (in this case, right and bottom sides of the left image). In addition style transfer on padded images has at least subjectively more pleasant colors, less artifacts and keeps a bit more content information and detail. The style used was the Edvard Munch's Scream painting.

273 be more reminiscent of blobs of colors and random  
274 artifacts, rather than the original content.

### 275 **7.1 Subjective aesthetic evaluation**

276 We compared our models that give best results with pre-  
277 trained flow models with style transfer modules being  
278 AdaIn and WCT. Results of such comparison can be  
279 seen in the figure 8. Comparing our two models in the  
280 first 2 columns, the model with 2 flow blocks transfers  
281 more accurate colours. Also, model with the 1 flow  
282 block preserves better details but performs worse in  
283 transferring style in macro scale. Therefore, model with  
284 two flow blocks has overall more aesthetically pleasant  
285 results. AdaIn and wct have comparable results to our  
286 models. However in some cases (like the "city" or  
287 "cat" image) AdaIn and WCT do not have as accurate  
288 colours taken from style image. Both preserve some  
289 colours from original content image, which if you want  
290 to do complete style transfer can be seen as a worse  
291 result. In addition, even though our proposed models  
292 should be biased by a style a bit, they preserves more  
293 accurate details from content image than the AdaIn.  
294 For example in the "city" image (content) and "pencil  
295 drawing" image (style), our models have more accurate  
296 reflections on the water.

### 297 **7.2 Multi-Style transfer experiments**

298 To verify if the value of attention based style trans-  
299 fer module was preserved our proposed model, we  
300 tried experimenting with interpolation of two different  
301 styles and use such new styles for the style transfer on  
302 the content image.

303 We took two images representing styles, extract  
304 z-vector from each style image in forward pass of  
305 neural network. Then interpolate these two vectors  
306 creating multiple new style z-vectors, where new z-  
307 vectors applied on the content show both styles with

308 more or less biases towards one of the styles. This  
309 shows run-time flexibility of the proposed model, that  
310 can be used for more interesting experiments without  
311 need of training it from scratch, e.g.: for tasks similar  
312 to multi-style transfer. The results of this experiment  
313 can be seen in the figure 5.

314 The architecture with only one flow block has  
315 visibly worse coloring than the one with two flow  
316 blocks, however, the interpolation between the two  
317 styles works either way.

### 318 **7.3 Reverse style transfer**

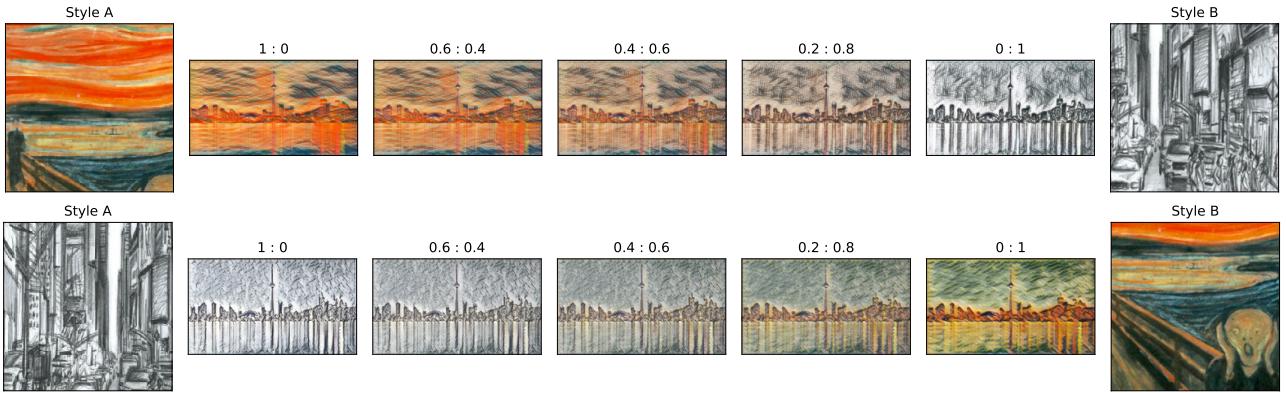
319 As discussed in section 3.4, the AdaAttN style transfer  
320 module is not unbiased. Therefore, it is highly proba-  
321 ble that some content information will be lost during  
322 the transfer process.

323 This can easily be demonstrated by styling one  
324 source image with an arbitrary style and then trying to  
325 replace this new style with the style of the original con-  
326 tent image. ArtFlow in combination with an unbiased  
327 style transfer module such as AdaIN or WCT will be  
328 able to perfectly reconstruct the original content image.  
329 However, since AdaAttN is not unbiased, ArtFlow in  
330 combination with this particular style transfer module  
331 will not.

332 This reverse style transfer capabilty can be ob-  
333 served in figure 6.

### 334 **7.4 Content preservation via ImageNet classi- 335 fication**

336 Another way to test how well the content information  
337 is preserved in the stylized images is to try to classify  
338 the objects in the images via a pre-trained network.  
339 For this purpose, we used an ImageNet [2] sample  
340 of 5000 images – 1000 classes of five examples each.  
341 This dataset was then stylized using first the original  
342 ArtFlow AdaIN model and then with our two AdaAttN  
343



**Figure 5.** Interpolation of two style (far left and far right) for two different network architectures. On the top, the network has 2 flow blocks each with 8 flow steps. The bottom image was generated by an architecture consisting of 1 flow block with 16 flow steps. Numbers above images represent proportion of given style images in applied style.



**Figure 6.** Reverse style transfer demonstration. The top picture is the content, second is the style, third is the stylized image and the last image is the reconstructed original content image. On the left, ArtFlow with AdaIN is perfectly capable of replacing the new style with the original content image style. On the right, the bias of AdaAttN makes it impossible to do so – the reconstructed original image retains some of the style characteristics, some detail is gone and the image is a bit distorted as a whole.

343 ArtFlow variants, one with 2x8 and the other with 1x16  
344 flow steps.

345 We then use three models that were pre-trained on  
346 the ImageNet dataset to classify these stylized images.  
347 The results can be seen in figure 7.

Once again, it is apparent that the unbiased nature 348  
of AdaIN allows for better preservation of the content 349  
factor. Interestingly enough though, the classification 350  
accuracy for our 2x8 AdaAttN network is a bit higher 351  
than for the 1x16 network, even though we would gen- 352  
erally consider the 2x8 model to be more 'aggressive' 353  
in the styling. 354

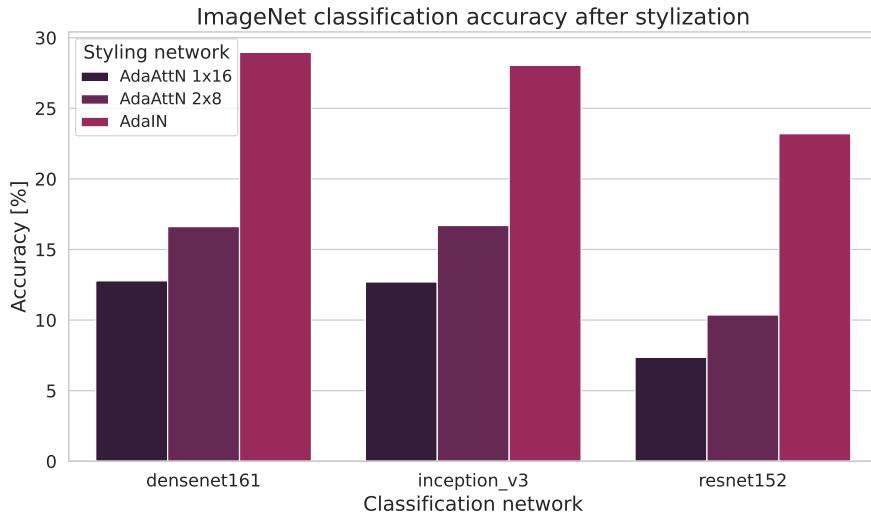
## 8. Conclusions

In this project, we combined two novel approaches to 355  
image style transfer, ArtFlow and AdaAttN in hopes 357  
of retaining the content clarity provided by ArtFlow 358  
and the stylistic aspects provided by the AdaAttN style 359  
transfer module. The resulting trained systems are 360  
producing decent results, held back only by some grid- 361  
shaped artifacts caused by the squeeze layers of the 362  
ArtFlow neural network. We evaluated the resulting 363  
systems both in terms of subjective aesthetic analysis 364  
and also by performing some objective experiments, 365  
namely interpolating between different styles and try- 366  
ing to classify the content of stylized images. 367

Computational resources were supplied by the project 368  
"e-Infrastruktura CZ" (e-INFRA CZ LM2018140 ) sup- 369  
ported by the Ministry of Education, Youth and Sports 370  
of the Czech Republic. 371

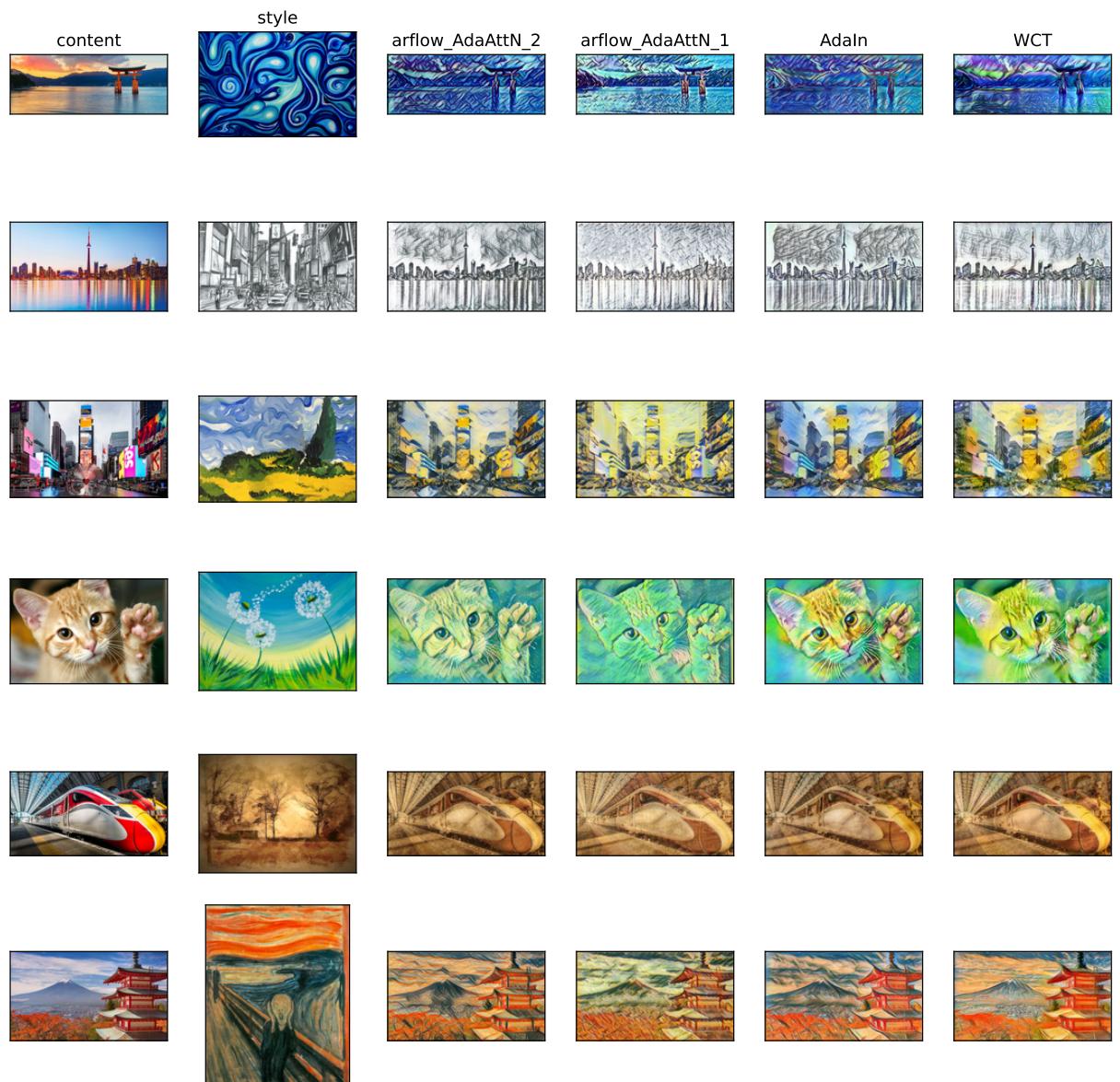
## References

- [1] AN, J., HUANG, S., SONG, Y., DOU, D., LIU, W. et al. ArtFlow: Unbiased image style transfer via reversible neural flows. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021. 373  
374  
375  
376  
377
- [2] DENG, J., DONG, W., SOCHER, R., LI, L.-J., LI, K. et al. Imagenet: A large-scale hierarchical im- 378  
age database. In: *Ieee. 2009 IEEE conference on* 379  
380



**Figure 7.** ImageNet classification results after stylizing the data.

- 381      computer vision and pattern recognition. 2009,  
 382      p. 248–255.
- 383 [3] GATYS, L. A., ECKER, A. S. and BETHGE, M.  
 384      *A Neural Algorithm of Artistic Style*. arXiv,  
 385      2015. Available at: <https://arxiv.org/abs/1508.06576>.
- 386 [4] HUANG, X. and BELONGIE, S. *Arbitrary Style  
 387 Transfer in Real-time with Adaptive Instance Nor-  
 388 malization*. arXiv, 2017. Available at: <https://arxiv.org/abs/1703.06868>.
- 389 [5] KINGMA, D. P. and DHARIWAL, P. *Glow: Gen-  
 390 erative Flow with Invertible 1x1 Convolutions*.  
 arXiv, 2018. Available at: <https://arxiv.org/abs/1807.03039>.
- 391 [6] LI, Y., FANG, C., YANG, J., WANG, Z., LU,  
 392      X. et al. Universal Style Transfer via Feature  
 393      Transforms. In: GUYON, I., LUXBURG, U. V.,  
 394      BENGIO, S., WALLACH, H., FERGUS, R. et al.,  
 395      ed. *Advances in Neural Information Process-  
 396 ing Systems*. Curran Associates, Inc., 2017.  
 397      Available at: <https://proceedings.neurips.cc/paper/2017/file/49182f81e6a13cf5eaa496d51fea6406-Paper.pdf>.
- 398 [7] LIN, T.-Y., MAIRE, M., BELONGIE, S., BOUR-  
 399      DEV, L., GIRSHICK, R. et al. *Microsoft COCO:  
 400      Common Objects in Context*. arXiv, 2014. Avail-  
 401      able at: <https://arxiv.org/abs/1405.0312>.
- 402 [8] LIU, S., LIN, T., HE, D., LI, F., WANG, M. et al.  
 403      AdaAttN: Revisit Attention Mechanism in Arbi-  
 404      trary Neural Style Transfer. In: *Proceedings of  
 405      the IEEE International Conference on Computer  
 406      Vision*. 2021.
- [9] LUGMAYR, A., DANELLJAN, M., VAN GOOL, L. and TIMOFTE, R. *SRFlow: Learning the Super-Resolution Space with Normalizing Flow*. arXiv, 2020. Available at: <https://arxiv.org/abs/2006.14200>.
- [10] SHENG, L., LIN, Z., SHAO, J. and WANG, X. Avatar-Net: Multi-scale Zero-shot Style Transfer by Feature Decoration. In: *Computer Vision and Pattern Recognition (CVPR), 2018 IEEE Conference on*. 2018, p. 1–9.



**Figure 8.** Comparison of proposed models on sub-sample of content and style images.