
The Final Project "Identification and Synthesizing of Raphaels paintings from the forgeries"

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Artistic identification and image style transfer has long been considered as a
2 challenging task due to the complex and unstructured nature of fine arts, especially
3 painting. In this report, we explore algorithms for identification and synthesizing
4 of Raphael's paintings from the forgeries. For identification problem, we first
5 use Geometric tight frame and Gabor wavelet to extract features. With selected
6 feature based on area under ROC curve(AUC), we apply several techniques such as
7 SVM, Neural-Network and Decision-Tree for classification. Leave-one-out cross
8 validation is used to ensure accuracy. As for synthesization problem, we follow
9 the method proposed in [2]. With VGG network and well-defined content and
10 style loss functions, style transfer is reduced to an optimization problem. The
11 highlights of our report is that not only detailed discussion of these methods are
12 presented, but we also endeavor to implement our own original thoughts to improve
13 these models.

14 1 Introduction

15 The project is divided as two parts. The first part is a classification problem based on 28 Raphael
16 paintings which have been labeled as authentic works and forgeries. We intend to use some feature
17 extraction techniques to first record the main characteristic of the images by some simple statistics,
18 and then apply the machine learning methods for classification problem such as SVM and Neural-
19 Network. The second part will be somewhat more complicated. We need to transfer the style of
20 Raphael paintings to another picture while maintaining the semantic content of it.

21 Task1

22 Distinguishing genuine paintings from fake ones is not easy task because forgeries tend to imitate
23 characteristic of the genuine ones. Algorithms such as Neural Network, SVM are quite popular for
24 classification. However, simply apply these algorithms seems like a bad idea because our data set,
25 usually over 10 thousand pixels one image, can be computational costly. Besides, as is mentioned in
26 [2], in higher layers of the network, detailed pixel information is lost, which is key information for
27 art authentication.

28 As a result, feature extraction without loss of information containing particular traces of style is an
29 important data preparation. The article [9] use 18 filters named geometric tight frame to convolve the
30 image before further experiments. Gabor wavelet filters, mentioned in is another efficient tool for
31 feature extraction. Leave-one-out cross validation shows our best performance reaches around 90%
32 accuracy with Gabor wavelet based decision tree model.

33 **Task2**

34 So far, many algorithms has been designed to accomplish the second task. The article [8] is based
35 on a bidirectional analogy of two semantically similar pictures, which is an important point since
36 the two picture will have very similar content information at the coarsest layer of the VGG network.
37 And the desired picture at this layer is initialized the same as the input. Then a nearest neighbor field
38 search method helps to reconstruct the picture in the previous CNN layer.
39 Another approach described in [2] to address the problem is more straightforward. The content and
40 style content are separately dealt with and a loss function is defined as a weighted sum of the ones
41 from the content and the style based on the VGG network. Then apply gradient descent method to
42 find the output image which minimize the loss function.

43 **2 Related Work**

44 **Task 1**

45 Li et al. [7] made an effort to extract those visually salient brushstrokes of van Gogh based on
46 an integrative technique of both edge detection and clustering based segmentation. With the ex-
47 tracted brushstrokes, some definitions of brushstroke features for art authentication were given in
48 distinguishing van Gogh paintings from forgeries.

49 Qi. et al. [11] use background selection and wavelet-HMT-based Fisher information distance for
50 authorship and dating of impressionist and post-impressionist paintings. Two novel points were
51 introduced in this work. The first point is that background information is much more reliable than
52 the details of an intricate object which cannot represent the artist's natural style because of multiple
53 edits and corrections. The second point is that an artist's style should be interpreted as a probability
54 distribution over a set of possible textures, and not just simply from the textures themselves.

55 **Task 2**

56 *Liao et al* [8] propose an attribute transfer technique across images that may have different appearance
57 but have perceptually similar semantic structure. They use a Convolutional Neural Network to
58 construct a feature space in which to form analogies subject to bidirectional constraints. What's more,
59 they adapt patch match and nearest neighbor field strategy to reconstruct the combined image.

60 *Gatys et al* [2] show how the generic feature representations learned by Convolutional Neural
61 Networks can be used to independently process the content and style of natural images it is a
62 texture transfer algorithm that constrains a texture synthesis method by feature representations from
63 Convolutional Neural Networks. Since the texture model is also based on deep image representations,
64 the style transfer method elegantly reduces to an optimization problem within a single neural network.
65 New images are generated by performing a pre-image search to match feature representations of
66 example images.

67 *Isola et al* [4] introduce a general solution to image transfer problem by training the network to not
68 only learn the feature map but also a loss function to establish this map automatically. The goal is
69 accomplished by Generative Adversarial Networks learning an appropriate loss function for a specific
70 task.

71 **3 Data Description**

72 The data set provided by Prof.Yang Wang from HKUST consists of high resolution scans of 28
73 paintings, scaled to a uniform density of 200 dots per painted inch. The picture size vary from
74 1192*748 to 6326*4457 pixels. Of the 28 paintings, 12 have been attributed to Raphael, 9 are known
75 to be non-Raphael, and others are currently questioned by experts.

76 Figure 1 is taken from our data set.

77 In section 4, for easier studying and feature capturing, we change them into grey-scale scans. Usually,
78 turning a colored picture into grey-scale might lose a lot of information. In this specific problem
79 however, the pictures given are all single-colored. So grey-scale doesn't make any difference.



(a) No.8 Raphael



(b) No.11 non Raphael

Figure 1: two pictures from data set

80 There are numerous ways for transformation. One simple and widely used formula is from BT.601, a
81 standard issued in 1982.

$$\text{Grey} = \text{Red} * 0.299 + \text{Green} * 0.587 + \text{Blue} * 0.114 \quad (1)$$

82 It's worth noting that our data set is quite small, only 21 labeled pictures. For synthesization this is
83 ok but not enough for classification. In an attempt to enlarge this data set, we cut each labeled image
84 into 4 parts. Thus we get 84 labeled data in total.

85 4 Artistic Authentication of Raphael's Paintings

86 4.1 Feature Engineering

87 Artistic authentication requires delicate analysis of specific style and habits. Although perceptual
88 meaningfulness is a subjective concept, it can be measured with a computer by using techniques
89 derived from biological and psychological models of the human visual system. The key point is to
90 find the appropriate features which give a good separation between the artists' paintings and those by
91 his imitators. Next we will give thorough illustration of two analysis tools, namely, Geometric tight
92 frame in [9] and Gabor wavelet transformation in [5].

93 Geometric Tight Frame Filters & Moment Statistics

94 The Geometric tight frame we use has 18 filters. They can capture the first and second-order
95 differences in the horizontal, vertical and diagonal directions in every small neighborhood of the
96 paintings. The first 2 filters are shown below. Refer to [9] to check all the filters.

$$\tau_0 = \frac{1}{16} \times \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}, \tau_1 = \frac{1}{16} \times \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} \quad (2)$$

97 Given the i -th color painting, with m_i -by- n_i pixels, we represent its grey-scale intensity by an
98 m_i -by- n_i matrix \mathbf{P}_i . Convolve \mathbf{P}_i with each τ_j , $0 \leq j \leq 17$, to get the corresponding m_i -by- n_i
99 tight frame coefficient matrices $\mathbf{A}^{(i,j)}$:

$$\mathbf{A}^{(i,j)} = \mathbf{P}_i * \tau_j = \begin{pmatrix} a_{1,1}^{(i,j)} & \dots & a_{1,n_i}^{(i,j)} \\ \vdots & & \\ a_{m_i,1}^{(i,j)} & \dots & a_{m_i,n_i}^{(i,j)} \end{pmatrix} \quad (3)$$

100 Therefore, there are 18 corresponding coefficient matrices for each painting after the decomposition. In
101 an attempt to reduce dimensionality of data set, [9] proposes the following 3 statistics as features.

(i) the mean of the entries in the coefficient matrix $\mathbf{A}^{(i,j)}$:

$$\mu^{(i,j)} = \frac{1}{m_i n_i} \sum_{l=1}^{m_i} \sum_{k=1}^{n_i} a_{l,k}^{(i,j)}$$

(ii) the standard deviation of the entries in $\mathbf{A}^{(i,j)}$:

$$\sigma^{(i,j)} = \left(\frac{1}{m_i n_i - 1} \sum_{l=1}^{m_i} \sum_{k=1}^{n_i} (a_{l,k}^{(i,j)} - \mu^{(i,j)})^2 \right)^{\frac{1}{2}}$$

(iii) the percentage of the "tail entries", which are entries that lies outside one standard deviation from the mean:

$$p^{(i,j)} = \#(\hat{\mathbf{A}}^{(i,j)}) / (m_i n_i).$$

where $\#(\hat{\mathbf{A}}^{(i,j)})$ is the number of nonzero entries in the tail matrix $\hat{\mathbf{A}}^{(i,j)}$ which is defined by:

$$\hat{a}_{l,k}^{(i,j)} = \begin{cases} a_{l,k}^{(i,j)}, & \text{if } |a_{l,k}^{(i,j)} - \mu^{(i,j)}| \geq \sigma^{(i,j)} \\ 0, & \text{otherwise} \end{cases}$$

Thus the feature vector of the i -th painting is represented by

$$[\mu^{(i,0)}, \dots, \mu^{(i,17)}, \sigma^{(i,0)}, \dots, \sigma^{(i,17)}, p^{(i,0)}, \dots, p^{(i,17)}] \in \mathbb{R}^{54}.$$

102 Gabor Wavelet Filters & Energy Values

Multiscale-oriented Gabor wavelet filters are sensitive to particular wave-like patterns. The Gabor wavelet filters come in pairs: $G_{even}(x, y, \sigma, \alpha, \omega)$ and $G_{odd}(x, y, \sigma, \alpha, \omega)$ are, respectively, the real and imaginary parts of the function:

$$e^{2\pi i \omega(x \sin \alpha + y \cos \alpha)} e^{-\frac{x^2+y^2}{2\sigma^2}}$$

103 where x and y are the spatial coordinates assuming the filters are centered at the origin, α sets the
104 spatial orientation and ω the spatial frequency of the filter pair.

105 In theory, there are infinite number of Gabor filters. In an attempt to balance efficiency and
106 convenience, we apply a set of filters that covers a range of orientations and scales: six orientations
107 ($\alpha = (k/6)\pi$, for $k \in \{0, 1, 2, 3, 4, 5\}$) and four scales (manually set ω and σ).

Figure 2 from [5] visualizes 24 filters.

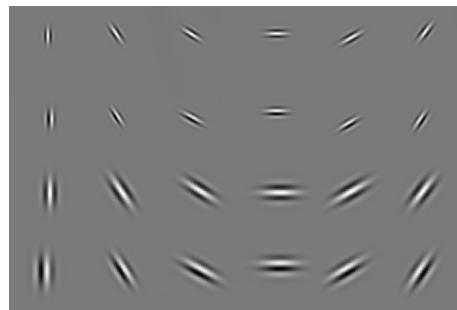


Figure 2: 24 Gabor filters

108

109 Convolving these filter with an image yields a decomposition of the image structure into "energy
110 values", one energy value for each pixel, orientation, and scale. The Gabor wavelet energy value of
111 an image patch is defined as the sum of the squared values obtained by convolving both components
112 with an image patch. Figure 3 from [5] is a good illustration of the whole process.

113 4.2 Model Implementation

114 We use the following models to do feature selection and classification.

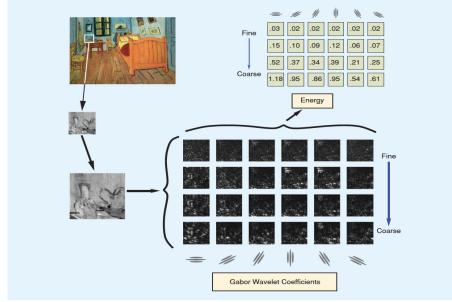


Figure 3: filters and energy value

115 **Forward Stage-Wise Feature Selection**

116 It is unlikely that all the 54 features obtained by Geometric frame are discriminative between Raphael
117 and his contemporaries. If we include some noisy features in our classification task, the accuracy will
118 deteriorate. That's why we adopt stage-wise method in [9].

119 To be precise, suppose there are 21(in this problem, we have 21 pictures labeled Raphael and
120 non-Raphael). let the training dataset be $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_{21}\}$ and

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{21} \end{bmatrix} \quad (4)$$

121 be the data matrix of \mathcal{X} and $\tilde{\mathbf{X}}$ be the normalization of \mathbf{X} such that each column in $\tilde{\mathbf{X}}$ has a unit
122 standard deviation. Let $\{1, \dots, 21\} = \mathcal{T}_{vG} \cup \mathcal{T}_{nvG}$ where \mathcal{T}_{vG} denotes the set of vG paintings in \mathcal{X} .
123 For any feature subset \mathcal{F} , denote $|\mathcal{F}|$ the number of elements in the set \mathcal{F} and $\mathcal{F} = \{i_1, \dots, i_{|\mathcal{F}|}\}$.
124 Define $\tilde{\mathbf{X}}_{j,\mathcal{F}} = (\tilde{\mathbf{X}}_{j,i_1}, \dots, \tilde{\mathbf{X}}_{j,i_{|\mathcal{F}|}})$. Then we define the vG center w.r.t. \mathcal{F} as the mean vector of vG
125 on \mathcal{F} , i.e.

$$\mathbf{c}^{\mathcal{F}} = \frac{1}{|\mathcal{T}_{vG}|} \sum_{j \in \mathcal{T}_{vG}} \tilde{\mathbf{X}}_{j,\mathcal{F}} \quad (5)$$

126 With this we define the distance between the j -th painting in \mathcal{X} and the vG center $\mathbf{c}^{\mathcal{F}}$ by

$$d_j^{\mathcal{F}} = \|\tilde{\mathbf{X}}_{j,\mathcal{F}} - \mathbf{c}^{\mathcal{F}}\|_2 \quad (6)$$

127 For \mathcal{F} to be a good feature set, $d_j^{\mathcal{F}}$ should be small for $j \in \mathcal{T}_{vG}$ and large for $j \in \mathcal{T}_{nvG}$ i.e. nvG
128 should be far from the vG center and regarded as outliers.

129 We apply the theory of the ROC curve. Sort $\{d_j^{\mathcal{F}}\}_{j=1}^{21}$ in an ascending order. $d_{j1}^{\mathcal{F}} \leq \dots \leq d_{j21}^{\mathcal{F}}$. For
130 a given number ρ , we can use it as a binary classifier and calculate the TPR and FPR, then draw a
131 ROC curve based on it. AUC denote area under the curve. Ideally, $AUC(\mathcal{F})=1$. We need to choose a
132 subset \mathcal{F} which maximize AUC, which can be hard to calculate. So we start from empty set $\mathcal{F}^{(0)}$ and
133 greedy choose the next feature by

$$l_{j+1} = argmax_{l \notin \mathcal{F}^{(j)}} AUC(\mathcal{F}^{(j)} \cup \{l\}) \quad (7)$$

134 In this way, we are able to select our labels both fastly and accurately.

135 **Support Vector Machine**

136 Since we have learned this algorithm this semester, we will only briefly summarize what we have
137 done. We use a linear model without kernel. Besides, we set the penalty coefficient equal 1.

138 **Decision Tree**

139 Python Machine Learning tool *scikit-learn* already has a library called *tree*. To implement this
140 algorithm, we simply need to implement a line of Python code:

141 `from sklearn import tree`
142

144 **Neutral Network**

145 We use 1-hiden-layer neutral network with 7 hidden units, the connection between layers are liner,
146 liner and softmax. In the Model fitting phase, we applied a dropout layer on the hidden layer(dropout
147 probability 0.5) to avoid over fitting and we use various step size. More specifically, the model can
148 automatically check the loss every 50 epochs, if the value not decrease for continuous 4 times, the
149 step-size will reduce by half.

150 **4.3 Experiments**

151 The following are essential codes to conduct the experiment, partitioned by their applications.

152 **4.3.1 Essential Codes**

153 **Packages to Incorporate**

```
154 import numpy as np
155 import pickle
156 import tensorflow as tf
157 import sklearn
158 from sklearn.svm import SVC
159 from sklearn import neighbors
160 from sklearn import tree
161 import time
162 from Feature_model import Tight_frame_classifier
```

165 **Moment Statistics and Energy Value**

```
166 def statistics(image):
167     """return the mean, std, precent of the tail and energy value"""
168     a1 = image.mean()
169     a2 = np.std(image)
170     n, m = image.shape
171     count = sum(sum(abs(image - a1) > a2))
172     a3 = count / (m * n)
173     a4 = sum(sum(image**2))
174     return a1, a2, a3, a4
```

177 **Gabor Wavelet**

```
178 def gabor_fn(sigma, alpha, omega):
179     """In our problem, sigma is std of Gaussian part,
180         alpha is orientation, omega is wave frequency"""
181     max = 30
182     min = -max
183     (y, x) = np.meshgrid(np.arange(min, max + 1), np.arange(min, max + 1))
184
185     # Rotation
186     x_alpha = x * np.cos(alpha) + y * np.sin(alpha)
187     y_alpha = -x * np.sin(alpha) + y * np.cos(alpha)
188
189     gb = np.exp(-.5 * (x ** 2 / sigma ** 2 + y ** 2 / sigma ** 2)) * np.cos(2 *
190         np.pi * omega * x_alpha)
191
192     return gb
```

194 **Stage-Wise-Forward Feature Selection**

```
195 def select_feature(train, label, size_feature, selected_feature=[]):
196     """
```

```

198     function:use stage wise forward algorithm to choose best
199     n-dimension features
200
201     input:Training matrix(list, 20*Feature_size), Training label,
202     size of selected feature (integer from 1 to Feature_size),
203
204     output:the selected feature(list with integer from 1 to Fe-
205     ature_size), corresponding AUC(0 to 1), std of training
206     matrix by column(np.array 1*Feature_size), mean vector of vG
207     (np.array 1*Feature_size), and threshold of distance(a real number).
208     """
209
210     Feature_size = train.shape[1]
211     # normalization
212     std_train = train.std(axis=0)
213     mean_train = train.mean(axis = 0)
214     norm_train = (train-mean_train) / std_train
215     norm_N = norm_train[label == 0, :]
216     norm_T = norm_train[label == 1, :]
217     center_true = norm_T.mean(axis=0) # mean vector of vG, i.e. vG center
218
219     auc_select = []
220     while (len(selected_feature)<=size_feature):
221         index_auc_pairs = []
222         unselected_feature = [i for i in range(Feature_size)
223                               if i not in selected_feature] # optional feature
224         for F in unselected_feature:
225             temp_feature = selected_feature + [F]
226             N_temp = norm_N[:, temp_feature]
227             T_temp = norm_T[:, temp_feature]
228             center_true_temp = center_true[temp_feature]
229             # distance between fitures and vG center L2-distance
230             dist_N = (((N_temp - center_true_temp) ** 2).sum(axis=1)) ** 0.5
231             dist_T = (((T_temp - center_true_temp) ** 2).sum(axis=1)) ** 0.5
232             dist = np.append(dist_T, dist_N)
233             label_temp = np.append(np.ones(T_temp.shape[0]),
234                                   np.zeros(N_temp.shape[0]))
235             # calculate and store AUC
236             fpr, tpr, thresholds = metrics.roc_curve(
237                 label_temp, dist, pos_label=0, drop_intermediate=False)
238             index_auc_pairs += [[F, metrics.auc(fpr, tpr)]]
239         # choose the best one feature
240         index_auc = sorted(index_auc_pairs, key=lambda x: x[1])
241         feature_best = index_auc[-1][0]
242         auc_select = index_auc[-1][1]
243         selected_feature += [feature_best] # added the best one into selected
244         features
245
246         # choose best threshold
247         N = norm_N[:, selected_feature]
248         T = norm_T[:, selected_feature]
249         center_t = center_true[selected_feature]
250         N_dis = (((N - center_t) ** 2).sum(axis=1)) ** 0.5
251         T_dis = (((T - center_t) ** 2).sum(axis=1)) ** 0.5
252         temp_label = np.append(np.ones(T.shape[0]), np.zeros(N.shape[0]))
253         dist = np.append(T_dis, N_dis)
254         accuracy = []
255         for j in range(len(dist)): # test difference threshold
256             thre = dist[j]
257             predict = dist.copy()
258             predict[dist > thre] = 0
259             predict[dist <= thre] = 1
260             accuracy += [accuracy_score(temp_label, predict)]
261         tempt = np.argsort(accuracy)
262         thres = np.mean(dist[tempt[-2:]]) # the best threshold
263     return selected_feature, auc_select, mean_train, std_train, center_true, thres

```

264 **Python Class for Tight Frame classifier**

```

265
266 class Tight_frame_classifier(object):
267     def __init__(self,best_size,Feature,mean_train,SD,Center,Threshold):\n
268         # store the best parameters of every possible dimensions of features\n
269         self.selected_feature = Feature\n
270         self.mean_train = mean_train\n
271         self.std_train = SD\n
272         self.center_true = Center\n
273         self.thres = Threshold\n
274         self.best_size = best_size\n
275\n
276     def best_feature(self):\n
277         return self.selected_feature[:self.best_size]\n
278\n
279     def predict(self,test,feature_size = None):\n
280         if (feature_size == None) :\n
281             feature_size = self.best_size\n
282         if (type(test[0]) == np.float):\n
283             test = [test]\n
284         prediction = []\n
285         sd = self.std_train[feature_size-1]\n
286         center = self.center_true[feature_size-1]\n
287         selected_feature = self.selected_feature[:feature_size]\n
288         thres = self.thres[feature_size-1]\n
289         for ite in test:\n
290             T = (ite-self.mean_train) / sd\n
291             T = T[selected_feature]\n
292             center_true_temp = center[selected_feature]\n
293             dist = (((T - center_true_temp) ** 2).sum()) ** 0.5\n
294             prediction += [dist < thres]\n
295     return prediction
296

```

297 **4.3.2 Results and Comparisons**

298 Remember that in order to enlarge our data set, we have cut each labeled image into 4 parts and get
299 84 labeled image. Then we randomly split our data into training set and test set (use Python sklearn
300 package). We adopt leave-one-out cross-validation on training set to select features, and tabular 1
301 shows our results.

Table 1: Experiment Results

Feature Extraction	model	TPR	TNR	Classification Accuracy
Tight Frame	Forward Stage wise	94.5%	34.7%	67.6%
	SVM	66.2%	80.9%	72.8%
	Decision Tree	75.6%	92.4%	83.2%
	KNN	47.3%	80.9%	62.4%
	Neural Network	76.2%	92.4%	84.0%
	Forward Stage wise	78.0%	78.0%	78.0%
Gabor Wavelet	SVM	100%	0	60%
	Decision Tree	95.3%	78.0%	88.4%
	KNN	95.3	0	57.2%
	Neural Network	95.3%	82.4%	90.1%

302 **Analysis**

303 The overall performance based on Gabor Wavelet seems better than Tight Frame filters, in that both
304 the best classification method(Decision tree model) and Forward Stage Wise method have higher
305 accuracy score. It's safe to conclude that Gabor Wavelet filters do a good job in feature extraction.

306 However, results of SVM & KNN based on Tight Frame features are better than that of Wavelet
307 features. One possible explanation is that both these two models' classification rule is based on
308 distance, while the selection of Geometric Tight features are also based on distance.

309 **Further Discussion: Gabor wavelet parameters**

310 Gabor wavelet is used in our report to extract of wave-like patterns. Different pictures may have
311 different scale, i.e. the reliability of our method largely depends on consistency of photography
312 techniques. Thanks to our specially made data set, with a uniform density of 200 dots per painted
313 inch, we don't need to worry about scale.

314 In our report, we manually check the scale of our data set and find that the stroke is around 5 to 20
315 pixels in width, ranging from the finest to coarsest. So we set a suitable parameter to ensure that all
316 the wave lengths roughly match the stroke width.

317 **5 Synthesization of Style Transferred Images by Raphael's Artworks**

318 The main idea for the task is to extract the low level features like the texture and color from one
319 image which is called the *style image* s , and apply it to some more semantic and high level features
320 on another *content image* c to derive the *style-transferred image* x . Roughly speaking, the *output*
321 *image* x should be the one preserving the objects and concepts in the *content image* c whilst deviating
322 little from the *style image* s in terms of the style.

323 We intend to practice the mechanism depicted in [2] instead of [8] because the style and content
324 images we want to extract feature from differ much both in their genre and semantic structure while
325 the latter is an important requirement for the image analogy manipulation in [8].

326 So intuitively, we can introduce a loss function to measure the distinction between x and s & c , all of
327 which are rearranged as column vectors. All we need to do then is to produce the *style-transferred*
328 *image* which minimizes the loss function. Consequently, the original style transfer objective reduces
329 to an optimization problem:

$$x^* = \underset{x}{\operatorname{argmin}} (\alpha \mathcal{L}_{content}(c, x) + \beta \mathcal{L}_{style}(s, x)) \quad (8)$$

330 where α and β are two hyper-parameters that allow us to control how much emphasis we want to
331 put on the content relative to the style, and $\mathcal{L}_{content}(c, x)$, $\mathcal{L}_{style}(s, x)$ are the content and style loss
332 functions separately.

333 Now the crucial insight in [2] is that the definitions of these content and style losses are based not on
334 pixel loss which is the Euclidean distance between the output x and the target c or s , but instead in
335 terms of some more perceptual differences between them. Such *semantic difference* are learnt by a
336 convolutional neural network.

337 **5.1 Feature Representations**

338 It turns out that convolutional neural networks pre-trained for image classification have already learnt
339 to encode perceptual and semantic information that we need to grasp the semantic differences. The
340 key point of [2] is that we can apply the VGGNet model pre-trained on the *ImageNet* data set so that
341 higher layers in the network capture the high level content in terms of objects and their arrangement
342 in the input image (Figure 4). The feature representations of style and content are separable and the
343 output image is synthesized by finding the one that simultaneously matches the content and the style
344 representation based on the loss functions defined by the activations of the network.

345 We shall first introduce some denotation for the structure of the VGG network. The network we
346 intend to work on is the 16-layer VGG network. Let $\phi^j(x)$ be the activations of the j -th layer of
347 the network ϕ when processing the image x . From [2], the content and style loss functions can be
348 defined as follows.

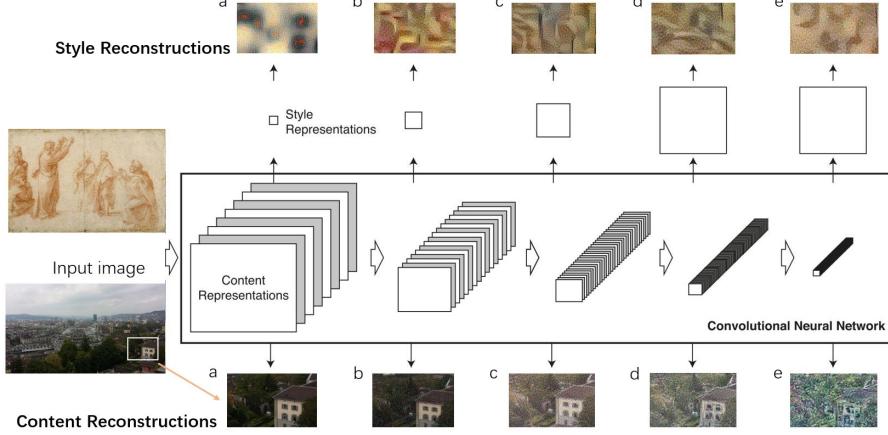


Figure 4: Image representation in a CNN. An input image is represented as a set of filtered images at each layer of CNN. The size of these filtered images is reduced by some down-sampling mechanism(e.g. pooling). To visualise the information in each layer, we reconstructed these filtered images.

349 5.2 Content Loss

350 For the content loss, the original choice in [2] is computed on the 4-th convolutional layer to be the
 351 squared Euclidean distance between feature representations of the content and combination images:

$$\mathcal{L}_{content}(\mathbf{c}, \mathbf{x}) = (\phi^j(\mathbf{c}) - \phi^j(\mathbf{x}))^2 \quad (9)$$

352 where $j = 4$. However, as stated above, the higher layer may deprive too much structure details of the
 353 input images and we will do experiments later to see which layer would be the most aesthetic choice.

354 5.3 Style Loss

355 The style representation will be more involved. According to [2], we first define a *Gram matrix*
 356 $G^l(x)$ which is proportional to the covariances of corresponding sets of features in the l -th layer. The
 357 elements of $G^l(x)$ is defined as:

$$G_{ij}^l(x) = \frac{1}{M^l} \sum_{k=1}^{M^l} \phi_{ik}^l(x) \phi_{jk}^l(x) \quad (10)$$

358 where ϕ_{ik}^l denotes the k -th position of the i -th vectorized feature map in the l -th layer and M^l is the
 359 number of elements in each feature map. The loss function of style representation in the l -th layer is
 360 then defined as the squared Frobenius norm of the difference between the Gram matrices of the style
 361 and combination images:

$$\mathcal{L}_{style}^l(\mathbf{s}, \mathbf{x}) = \|G^l(\mathbf{s}) - G^l(\mathbf{x})\|_F^2. \quad (11)$$

362 And the total style loss is a weighted sum of the contribution from each layer:

$$\mathcal{L}_{style}(\mathbf{s}, \mathbf{x}) = \sum_{l=1}^L w^l \mathcal{L}_{style}^l. \quad (12)$$

363 To see why the Gram matrix works, we refer to the paper [1]. It is said that due to the fact that the
 364 style is static and uniform across the image, the Gram matrix is used to capture global statistics by
 365 averaging over positions. Besides, coherence among multiple features (computed by different maps)
 366 needs to be preserved to capture information about which features tend to activate together, which is
 367 measured by the off-diagonal terms in the Gram matrix.

368 **5.4 Total Variation Regularization**

369 If we are going to solve the optimization problem with only the two loss terms defined above, the
370 output synthesized image would be quite noisy. According to the articles [6] and [10], we can add a
371 regularization term called the total variation loss to encourage spatial smoothness.

372 The total variation encourages images to consist of piece-wise constant patches. For continuous
373 functions $f : \mathbb{R}^2 \supset \Omega \rightarrow \mathbb{R}$, the total variation is given by:

$$\mathcal{R}(f, \delta) = \int_{\Omega} \left[\left(\frac{\partial f}{\partial u}(u, v) \right)^2 + \left(\frac{\partial f}{\partial v}(u, v) \right)^2 \right]^{\frac{\delta}{2}} dudv. \quad (13)$$

374 For input images $\mathbf{x} \in \mathbb{R}^{H \times W}$, the total variation is replaced by the finite discrete approximation:

$$\mathcal{R}(\mathbf{x}, \delta) = \sum_{i,j} ((x_{i,j+1} - x_{i,j})^2 + (x_{i+1,j} - x_{i,j})^2)^{\frac{\delta}{2}} \quad (14)$$

375 where δ is modified to balance the sharpness of the image against the removal of the artifacts called
376 "spikes". According to the paper *Fast image deconvolution using hyper-laplacian priors*, regularizer
377 with $\delta < 1$ is often used as a better match of the gradient statistics of natural images, while by
378 choosing $\delta > 1$ changes in images are distributed across regions rather than concentrated at a point
379 or curve.

380 In addition, when the target of the reconstruction is a colour image, the regulariser is summed for
381 each colour channel.

382 **5.5 Total Loss Function**

383 Based on the previous knowledge, the total loss function is then refined as:

$$\mathcal{L}(\mathbf{c}, \mathbf{s}, \mathbf{x}) = \alpha \mathcal{L}_{content}(\mathbf{c}, \mathbf{x}) + \beta \mathcal{L}_{style}(\mathbf{s}, \mathbf{x}) + \gamma \mathcal{R}(\mathbf{x}, \delta) \quad (15)$$

384 where α, β, γ and δ are four hyper-parameters we will check during the experiments. And the
385 objective is then to find the solution:

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} \mathcal{L}(\mathbf{c}, \mathbf{s}, \mathbf{x}). \quad (16)$$

386 **5.6 Experiments**

387 **5.6.1 Essential Codes**

388 The following are essential codes to conduct the experiment, partitioned by their applications.

389 **Packages to Incorporate**

```
390 import numpy as np
391 import time
392 import scipy
393 from PIL import Image
394 from os import listdir, mkdir, getcwd
395 from keras.applications.vgg16 import VGG16
396 from keras import backend as K
```

399 **Content Loss**

```
400 def content_loss(content, combination):
401     return K.sum(K.square(combination - content))
```

404 **Style Loss**

```
405 def gram_matrix(x):
406     features = K.batch_flatten(K.permute_dimensions(x, (2, 0, 1)))
407     gram = K.dot(features, K.transpose(features))
408     return gram
409
410 def style_loss(style, combination):
411     S = gram_matrix(style)
412     C = gram_matrix(combination)
413     channels = 3
414     size = height * width
415     return K.sum(K.square(S - C)) / (4. * (channels ** 2) * (size ** 2))
416
```

418 **Total Variation Loss**

```
419 def total_variation_loss(x):
420     a = K.square(x[:, :width-1, :height-1, :] - x[:, 1:, :height-1, :])
421     b = K.square(x[:, :width-1, :height-1, :] - x[:, :width-1, 1:, :])
422     return K.sum(K.pow(a + b, 1.5))
423
```

425 **Total Loss Function**

```
426 def total_loss(model, combination_image):
427     loss = K.variable(0.)
428
429     layers = dict([(layer.name, layer.output) for layer in model.layers])
430     layer_features = layers['block2_conv2']
431     content_image_features = layer_features[0, :, :, :]
432     combination_features = layer_features[2, :, :, :]
433
434     # content loss
435     loss += content_weight * content_loss(content_image_features,
436                                             combination_features)
437
438     # style loss
439     feature_layers = ['block1_conv2', 'block2_conv2',
440                       'block3_conv3', 'block4_conv3']
441
442     for layer_name in feature_layers:
443         layer_features = layers[layer_name]
444         style_features = layer_features[1, :, :, :]
445         combination_features = layer_features[2, :, :, :]
446         sl = style_loss(style_features, combination_features)
447         loss += (style_weight / len(feature_layers)) * sl
448
449     # total variation loss
450     loss += total_variation_weight * total_variation_loss(combination_image)
451
452     return loss # return total loss
```

454 **Optimization**

```
455 def eval_loss_and_grads(x):
456     # @x: the combination image.
457     x = x.reshape((1, width, height, 3))
458
459     outputs = [loss]
460     outputs += grads
461     f_outputs = K.function([combination_array], outputs)
```

```

463
464     outs = f_outputs([x])
465     loss_value = outs[0]
466     grad_values = outs[1].flatten().astype('float64')
467     return loss_value, grad_values
468
469 def minimize_loss(combination):
470     """
471         Using stochastic gradient descent, via fmin_l_bfgs_b algorithm.
472         @combination: concatenated name of the two images used.
473     """
474     x = np.random.uniform(0, 255, (1, width, height, 3)) - 128
475     evaluator = Evaluator()
476
477     print("\n\nProcessing: " + combination)
478     for i in range(iterations):
479
480         # print diagnostic information
481         print('Start of iteration', i)
482         start_time = time.time()
483         x, min_val, info = scipy.optimize.fmin_l_bfgs_b(evaluator.loss,
484                                         x.flatten(),
485                                         fprime=evaluator.grads,
486                                         maxfun=20)
487         print('Current loss value:', min_val)
488         end_time = time.time()
489         print('Iteration %d completed in %ds' % (i, end_time - start_time))
490     return x

```

492 We will perform experiments on two sets of images: one is the works of Raphael as style images and
493 the other is the test images, mostly landscape photographs served as content sources. First review the
architectures of VGG networks. Since we are not interested in the classification problem, we don't

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Figure 5: VGG Network Architectures

494 need the fully-connected layers or the final softmax classifier. We only need the part of the model
495 marked green in the table.

496 Besides, when decoding the original style image into a vector, we find that the artworks are stored in 4
497 channels, the last of which is the alpha compositing that represents the transparency. We will only
498 use the first 3 channels to make the dimension consistent with the content image. Also, to speed up
499 the reconstruction procedure, the original content image is condensed to a smaller size to relieve the
500 computation burden. Moreover, if the height of style image is larger than its width, we will rotate it
501 to fit the size of the content image.

502 The article [2] synthesizes the images by matching the content representation on layer 'conv4_2'
503 and the style representation on layers 'conv1_1', 'conv2_1', 'conv3_1', 'conv4_1' and 'conv5_1'.
504 We now implement this model with settings $\alpha/\beta = 1 \times 10^{-3}$, $\gamma = 1$ and $\delta = 1$ to have a glance of
505 the combined output. The combination image x begins its life as a white noise one and we use the
506 L-BFGS algorithm to iteratively improve upon it.

507 When we use the Raphael's artwork numbered 2 as the style image, the combination image is
produced in figure 6.



Figure 6: Images that combine the content of a photograph with the style of Raphael's artworks. The images were created by finding an image that simultaneously matches the content representation of the photograph and the style representation of the artwork. The ratio $\alpha/\beta = 10^{-3}$. The left one is the style image and the right one is the combination image.

508

509 Trade-off between Content and Style Representations When synthesizing an image that combines the
510 content of one image with the style of another, there usually does not exist an image that perfectly
511 matches both of them. However, since the loss function we minimize during image synthesis is a
512 linear combination between the loss functions for content and style respectively, we can smoothly
513 regulate the emphasis on either reconstructing the content or the style [2]. A high ratio α/β will
514 result in images that effectively capture the texture of the style artwork but showing hardly any of the
515 photograph's content.

516 Figure 7 shows the increasingly explicit content outline in the combination image by using larger
517 ratios. The combination image is produced after 10 iterations of the algorithm. It is apparent that
518 the output with ratio $\alpha/\beta = 10^{-2}$ produces the most aesthetic one that simultaneously maintains the
519 content and transfers the style based on the other hyper-parameters set previously.

520 In these results it is clear that the pre-trained network is aware of the semantic content of images,
521 especially objects like the parking cars. This is partly because the VGG networks are basically trained
522 on the classification data set comprised of particular elements. So the style transfer will preserve
523 these features more than the background [6].

524 **5.6.2 Effect of Different Layers for Content and Style Representation**

525 Another consideration in the image synthesizing process is the choice of layers for the extraction
526 of the style and content. According to [2], using the content representation in higher layers in the

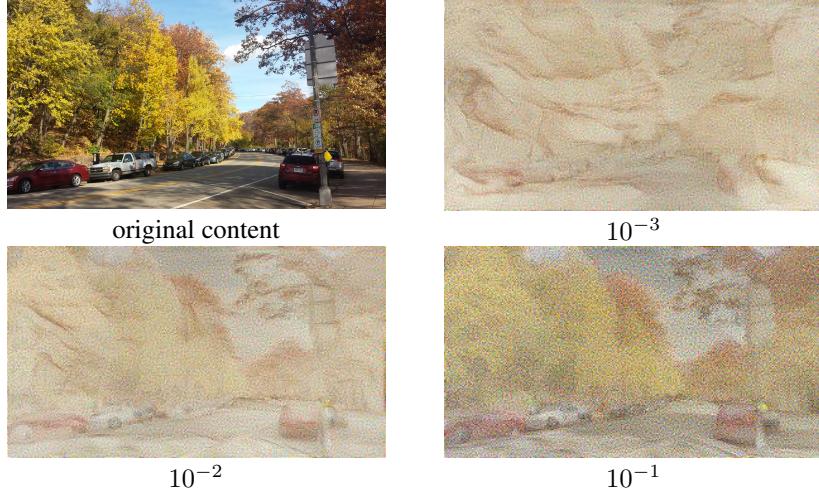


Figure 7: Output images synthesized by different settings of the ratio α/β . A high emphasis on the style produces combination difficult to recognize the content but excessively rich in the target style. A best choice for the model seems to be $\alpha/\beta = 10^{-2}$.

527 network decreases the object detail while matching the style representations up to higher layers
 528 preserves local textures within an increasingly large scale.

529 The previous experiment is conducted on the layer 'conv4_2' for content representation and 'conv1_1',
 530 'conv2_1', 'conv3_1', 'conv4_1' and 'conv5_1' for style representation. Alternatively, we perform
 531 the style transfer with the same artwork and parameter configuration ($\alpha/\beta = 10^{-2}$) but matching
 532 the photograph on layer 'conv2_2'. Figure 8 demonstrates the comparison between this kind of
 difference.

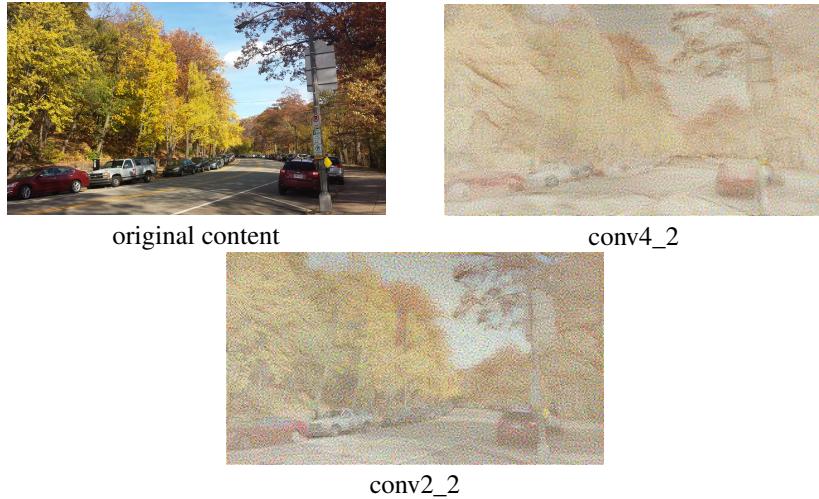


Figure 8: Matching the content on layer 'conv2_2' produces the combination image with more details but seems to hardly capture the style where the ratio $\alpha/\beta = 10^{-2}$.

533
 534 A more aesthetic choice by *Johnson et al.* is that we further replace the layer for style representation
 535 by 'conv1_2', 'conv2_2', 'conv3_3' and 'conv4_3'. The result better expresses the texture of the
 536 style image since we are performing the style representation on deeper layers which incorporate
 537 information from an increasingly larger spatial extent. We also change the ratio α/β expecting for a
 538 better balance between the content and style (Figure 9).

539 However, it looks like the style over the combination image is still somewhat beyond recognition.
 540 One explanation is that the style and content image differ greatly in their characteristics. For example,

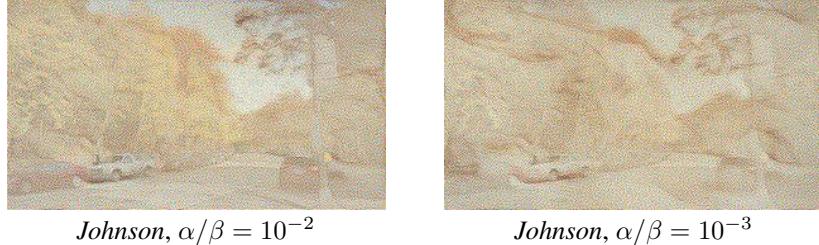


Figure 9: Reconstruction of the combination image with alternative layers for the style representation.

541 one is a sketch painting with sober beige color while the other is a colorful real-life photograph so that
 542 the texture and color from the style image may get blended with the content during the reconstruction
 543 process.

544 5.6.3 Selection of the Power Parameter δ of the Total Variation Regularization

545 To see how the total variation regularization term affects the combination image, we first draw a
 546 simple comparison between the output with and without this term. We will again use the style
 547 transfer structure performed in [6] and the ratio $\alpha/\beta = 10^{-2}$ and the power parameter $\delta = 1$. The
 548 regularization coefficient γ is set to 1. Figure 10 illustrates of the total variation regularization effect
 549 while the result seems to indicate the presence of the regularization term is out of usage. However, as
 550 we gradually increment the value of δ , the power parameter, the dramatic improvement will impress
 us.

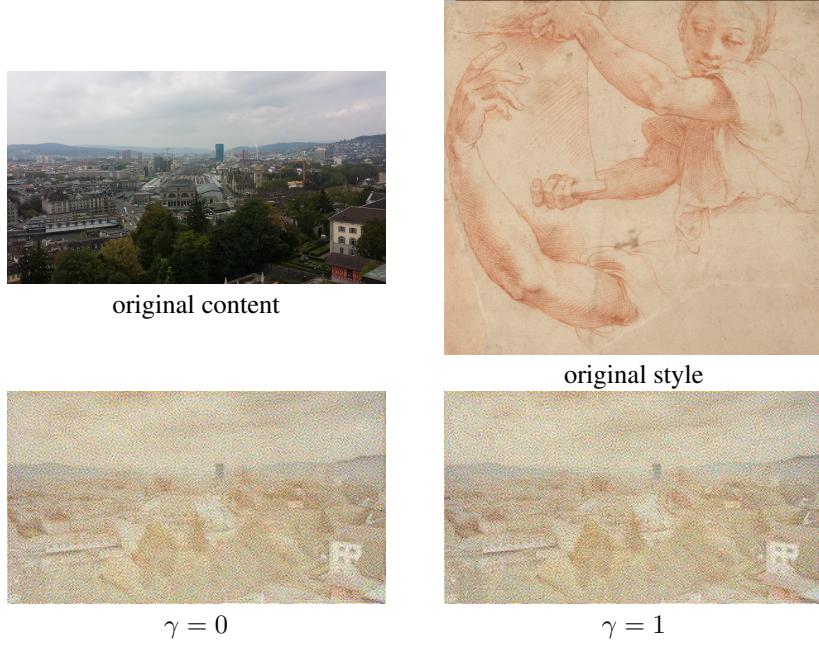


Figure 10: Reconstruction of style transferred image with and without total regularization when $\delta = 1$ shows the regularization produces little effect on removal of the 'spikes' which are sharp fluctuation of RGB value causing a discrete impression of image.

551
 552 Figure 11 shows the effect of increasing value of δ . It is obvious that the combination image
 553 becomes smoother and abnormal background transients disappear with larger values of δ , exactly
 554 what we desire. Future experiment and discussion can be done to find the optimal option for γ , the
 555 regularization coefficient. [10] presents us with a concrete theory of method to implement. Also,
 556 the regularization term can be defined as the α -norm $\mathcal{R}(\mathbf{x}) = \|\mathbf{x}\|_\alpha^\alpha$. By choosing a relatively large
 557 exponent the range of the image is encouraged to stay within a target interval instead of diverging.

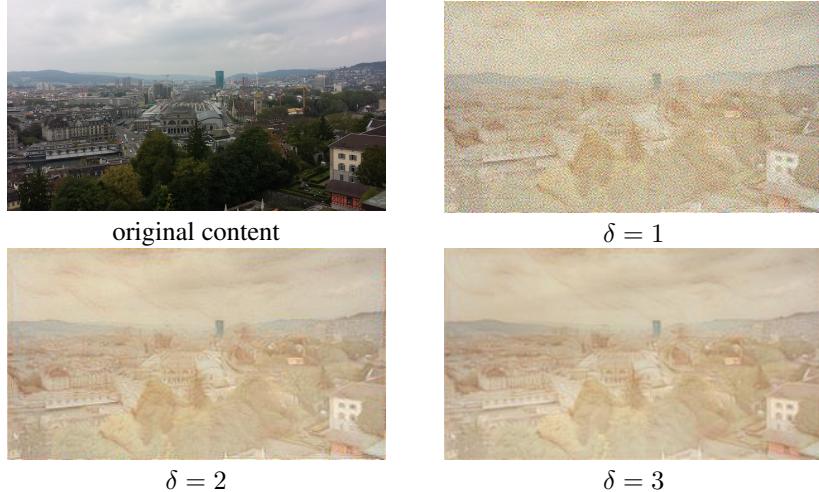


Figure 11: Increasing value of δ can bring about more spatial smoothness.

558 However, due to the limit of our personal laptops' computation ability (without GPU), we can only
 559 present a low-resolution combination image within a reasonable amount of time (i.e. one iteration of
 560 the output with 300 pixels of width costs over 10 minutes). So future improvement of this project
 561 will be focusing on finding a more efficient and time-saving algorithm, for example to build the style
 562 and content representation on the network architecture described in [12] and [3] and the feed-forward
 563 network in [6].

564 6 Conclusion

565 In task1, we use Tight Frame and Wavelet transform method to extract brushwork features, then
 566 use leave-one-out cross validation and stage-wise-forward algorithm to select features. After that,
 567 several models are used. From experiment result, it's obvious different feature extraction method
 568 significantly influence the efficiency of model and the dimension of features matters as well.

569 Of all models we trained in previous part, we did classification on the 7 pictures remained disputed
 570 (Pic 1/7/10/20/23/25/26). Basd on both *Tight Frame & Wavelet* features, most experiments show that
 571 Picture 1/7/23 are genuine, and Picture 10/25/26 are counterfeit. For Picture 23, predict results are
 572 various, so we have reservation about it.

573 In task2, to generate images with mixed content and style representations from two unrelated images,
 574 we used well-trained CNN network and *Gram matrix* to calculate the representations of content and
 575 style. We generate a new image from one noisy picture which minimize the pre-defined loss function.
 576 To enhance the effects of this model, we discussed different influence factors, including trade off
 577 between *style* and *content*, influence of different CNN layer etc.

578 References

- 579 [1] G. Berger and R. Memisevic. Incorporating long-range consistency in cnn-based texture
 580 generation. *CoRR*, abs/1606.01286, 2016.
- 581 [2] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *CoRR*,
 582 abs/1508.06576, 2015.
- 583 [3] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CoRR*,
 584 abs/1512.03385, 2015.
- 585 [4] P. Isola, J. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional
 586 adversarial networks. *CoRR*, abs/1611.07004, 2016.

- 587 [5] C. R. Johnson, E. Hendriks, I. Berezhnoy, E. Brevdo, S. Hughes, I. Daubechies, J. Li, E. Postma,
588 and J. Z. Wang. Image processing for artist identification – computerized analysis of vincent
589 van gogh’s painting brushstrokes. *IEEE Signal Processing Magazine*, July 2008.
- 590 [6] J. Johnson, A. Alahi, and F. Li. Perceptual losses for real-time style transfer and super-resolution.
591 *CoRR*, abs/1603.08155, 2016.
- 592 [7] J. Li, L. Yao, E. Hendriks, and J. Z. Wang. Rhythmic brushstrokes distinguish van gogh from
593 his contemporaries: Findings via automated brushstroke extraction. *IEEE Trans. Pattern Anal.*
594 *Mach. Intell.*, 34(6):1159–1176, 2012.
- 595 [8] J. Liao, Y. Yao, L. Yuan, G. Hua, and S. B. Kang. Visual attribute transfer through deep image
596 analogy. *CoRR*, abs/1705.01088, 2017.
- 597 [9] H. Liu, R. H. Chan, and Y. Yao. Geometric tight frame based stylometry for art authentication
598 of van gogh paintings. *CoRR*, abs/1407.0439, 2014.
- 599 [10] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them.
600 *CoRR*, abs/1412.0035, 2014.
- 601 [11] H. Qi, A. Taeb, and S. M. Hughes. Visual stylometry using background selection and wavelet-
602 hmt-based fisher information distances for attribution and dating of impressionist paintings.
603 *Signal Processing*, 93(3):541–553, 2013.
- 604 [12] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional
605 generative adversarial networks. *CoRR*, abs/1511.06434, 2015.