

# 画像スタイル変換と Web 画像を用いた画像の任意質感生成

松尾 真\*1 下田 和\*2 柳井 啓司\*3  
Shin Matsuo Wataru Shimoda Keiji Yanai

\*1 \*2 \*3電気通信大学大学院 情報理工学研究科  
Department of Informatics, The University of Electro-Communications

実世界画像の質感の任意変換は質感ベースの画像データの増量やデザイン、エンターテインメントなどの様々な分野での応用が期待できる。本研究では、2015 年に Gatys らが考案した、物体の形状を精密に維持して画像のスタイルを変換する Neural Style Transfer アルゴリズムを実世界画像の質感の転写に応用し、コンテンツ画像と質感単語による画像内物体の任意質感変換を目指す。

単語概念からの自動スタイル生成では、手動で行っていたスタイル画像の選択を自動化し、質感単語で収集した Web 画像から DCNN 中間層の出力と画像のスタイル認識のための特徴量 Neural Style Vector を用いてスタイルクラスタを構築し、スタイル表現を自動生成するシステムを構築する。

評価は生成画像へのアンケートを用いて行った。その結果、視覚的概念への連想が容易であり、類似スタイル表現を多く持つクラスタが生成される単語では人間に評価されやすい画像が生成され、またそのクオリティはコンテンツのカラーやスタイルとの共通構造によって、大きく影響を受けることが分かった。

## 1. はじめに

近年、2012 年の ImageNet Large Scale Visual Recognition Challenge(ILSVRC) のクラス分類タスクにおいて最も高い精度を挙げた Deep Convolutional Neural Network (DCNN) [3] がコンピュータビジョンの分野で注目されており、様々な分野へと応用する動きが進められている。

2015 年、Gatys ら [1] によって、Deep Neural Network(DNN) を用いた、画像のスタイルを絵画のものに変換するアルゴリズムが考案された。これにより、従来の画像合成技術よりも物体の形状を精密に維持して画像のスタイルを変換することが可能となった。

本研究では、スタイル変換アルゴリズムを Flickr Material Database(FMD) やオノマトペを用いて Web から収集した実世界画像の質感の転写に応用し、画像内物体の任意質感変換を目指す。

実世界画像の質感概念をスタイルとして、画像内物体の質感の任意変換が可能となれば、画像に対する心象を意図的に変化させることができ、画像データの増量やデザイン、エンターテインメントなどの様々な分野での応用が期待できる。

しかし、このアルゴリズムには目的のスタイルが画像の形式で存在する必要があるため、ある質感を与えたいときにはその質感を持つ画像をマニュアルで用意する必要がある。この問題の解決のため、画像のスタイル認識のために考案した特徴量 Neural Style Vector[2] を用いて、質感単語で収集した Web 画像から有用なものを自動選択し、スタイル表現を生成するシステムを構築する。

## 2. Neural Style Transfer

Deep Neural Network を用いたスタイル変換の関連研究は Gatys ら [1] の研究が挙げられる。この研究では、Imagenet1000 クラスを学習済みの DNN の中間層のフィルタ出力を使用することで、スタイル画像のスタイルをコンテンツ画像に転写しており、画像のクラス分類が中心であった DNN の研究分野に衝撃を与えた。

本研究では Gatys ら [1] の手法を用いて画像を合成することで、画像の質感の変換を行う。変換させる画像をコンテンツ画像  $x_c$ 、スタイル画像を  $x_s$ 、合成結果画像を  $x_g$  とする。

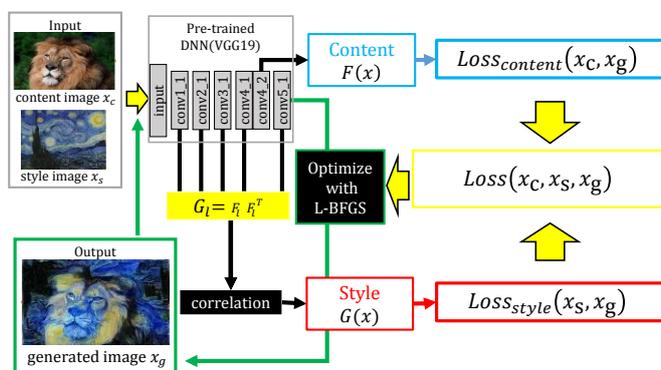


図 1: スタイル変換アルゴリズム

$x_c$ ,  $x_s$ ,  $x_g$  のコンテンツ表現とスタイル表現を CNN の特定の layer の活性化値から求め、 $x_g$  のコンテンツ表現が  $x_c$  に、スタイル表現が  $x_s$  に近くなるように反復的に合成する。

使用した CNN は VGG19[3] であり、コンテンツ表現に使用する layer は conv4\_2, スタイル表現に使用する layer は conv1\_2, conv2\_2, conv3\_4, conv4\_4, conv5\_4 である。図 1 にスタイル変換アルゴリズムの概略を記す。

layer  $l$  におけるコンテンツ表現はパラメータ数  $N_l$  の活性化値行列  $F_l(x)$ , その損失関数は  $x_c$  と  $x_g$  の差であり、式 1 で表される。

$$L_c(x_c, x_g) = \frac{1}{2} \sum_{i,j} (F_{l,i,j}(x_c) - F_{l,i,j}(x_g))^2 \quad (1)$$

layer  $l$  におけるスタイル表現は活性化値行列の式 2 で表される相関行列  $G(x, l)$ , その損失関数は  $x_s$  と  $x_g$  の差であり、式 3 で表される。使用する layer 全体の誤差は重み  $w_l$  を用いて式 4 で表される。

$$G_l(x) = F_l(x)F_l^T(x, l) \quad (2)$$

$$L_{loss_s, l}(x_s, x_g) = \frac{1}{4N_l^2} \sum_{i,j} (G_{l,i,j}(x_s) - G_{l,i,j}(x_g))^2 \quad (3)$$

$$Loss_s(x_s, x_g) = \sum_l w_l Loss_{s, l}(x_s, x_g) \quad (4)$$

全体のエラー関数は重み  $w_c, w_s$  を用いて式 5 で表される。この式の値が最小となるように  $x_g$  を L-BFGS 法を用いて最適化する。

$$Loss(x_c, x_s, x_g) = w_c Loss_c + w_s Loss_s \quad (5)$$

現在このアルゴリズムには様々な改良が考案されている。Li ら [4] は Markov Random Fields(MRF) と CNN の組み合わせによって実世界画像の転写により適したスタイル変換を行った。Novak ら [5] は Gram Matrix のシフトや layer 間の Gram Matrix の操作によってスタイル変換の質を向上させた。Berger ら [6] は特徴マップの空間転移により、スタイルを広範囲に保持したスタイル変換や Inpainting を行った。

また、Feed forward network を用いてスタイルを学習することによる高速化 (Fast style transfer) が Johnson ら [7] や Ulyanov ら [8] により考案された。さらに単一の Feed forward network で複数のスタイルを学習したマルチスタイルの Fast style transfer が Dumoulin ら [9] によって考案された。

### 3. 単語概念からの自動スタイル生成

類似スタイルを持つ画像同士が集中したクラスタをスタイル表現として構成するために、以下の手順で質感単語から収集された Web 画像から DCNN 特徴 (VGG19 の fc6 層の出力) と Neural Style Vector を用いて有用なスタイル表現の抽出を行い、そのスタイル表現によるスタイル変換の結果を評価する。また、質感単語をキーワードとした画像収集には Bing API を使用する。

1. 収集した画像を  $W_r \times W_r$  にリサイズ
2. 全ての画像から DCNN 特徴を抽出
3. 不適切な画像を除去
4. 残った画像から  $W_p \times W_p$  のパッチを切り出し
5. 全てのパッチから Neural Style Vector を抽出
6. DCNN 特徴を用いて画像をクラスターリング (クラスター数  $N$ )
7. Neural Style Vector を用いて各クラスター内の画像をリランキン
8. 各クラスター上位  $K$  枚の画像を用いてスタイル変換

実験では  $N = 10$  個のスタイルクラスターを構成し、各クラスターに式 6 でスコアを設定し、上位 3 個のクラスターを実験では使用した。クラスター  $C$  に含まれる画像数を  $Num(C)$ 、クラスター内の画像のクラスター中心との距離  $Dist(x_i, C)$  とする。

$$score(C) = \frac{N(C)}{\sum_i Dist(x_i, C) + \lambda} \quad (6)$$

実験では  $W_r = 1024, W_p = 512$  とした。収集した Web 画像から低解像度な画像や人間単体の画像を、DCNN 特徴で学習した SVM を用いて除外した。

異なる物体間で類似したスタイルを持つ例は少ないため、まず DCNN 特徴を用いて kmeans クラスターリングを行うことで、物体カテゴリ別のクラスターに分離させる。

各クラスターからのスタイル表現の導出のため、Neural Style Vector によるクラスター内のリランキン

グを行い、各クラスター内の上位  $K$  枚を選出する。リランキン

グにはまずクラスター内の画像をさらにクラスターリングし、各小クラスターに式 6 でスコアを与え、小クラスターをリランキン

グし、さらに各小クラスター内の画像をクラスター中心との距離が小さい順にリランキン

## 4. 実験

実験では FMD および質感単語で構築したスタイルクラスターを用いてスタイル変換を行い、ユーザー評価を行う。使用するスタイルクラスターは各単語の上位 3 つを選出し、クラスター上位 9 枚を使用してスタイル変換画像を生成する。

### 4.1 実験データ

コンテンツ画像は図 2 の 3 種類の画像を使用した。質感単語はオノマトペ 23 単語を使用した。



図 2: コンテンツ画像 (左から content0, 1, 2 とする)

### 4.2 ユーザー評価

ユーザー評価では図 3 のように、質感単語と正解の画像 3 枚、ランダムに選出した同コンテンツの不正解の画像 12 枚をブラウザに表示し、ユーザーが最もその単語に当てはまると感じた画像を選択し、その正解の可否で評価する。問題数はコンテンツ数 3, 単語数 23 語の 69 問、回答したユーザー数は 9 人だった。

### 4.3 評価結果

各単語のユーザー評価の結果は図 5 のようになった。正解率が高かった単語は「ふさふさ」、「ざらざら」、「もふもふ」だった。また、コンテンツ別の評価結果は図 6 のようになった。

各単語のユーザー評価の結果は図 5 のようになった。正解率が高かった単語は「ふさふさ」、「ざらざら」、「もふもふ」だった。正解率に繋がる要因として考えられるのは、「類似するスタイルクラスターの構成の成功」、「スタイル変換のクオリティ」、「質感単語の視覚的イメージのしやすさ」であると考えられる。

これらの単語のクラスターの内、最も選択されることが多かったものが図 4 であり、これらのスタイルクラスターは全て類似したスタイルが上位に集中しており、スタイル画像の収集が成功した例だと言える。

また、コンテンツ別の評価結果は図 6 のようになった。これらと比較すると質感単語の正解率はコンテンツによって大きく異なることが分かる。

図 7, 図 8, 図 9, 図 10 はそれぞれ「ふさふさ」、「ざらざら」、「もふもふ」、「くしゅくしゅ」から生成されたスタイルクラスターおよび変換画像とアンケートにおいての正解者の数である。

最も「ふさふさ」な画像を選んでください。0/72



図 3: ユーザー評価実験 (赤枠が正解)

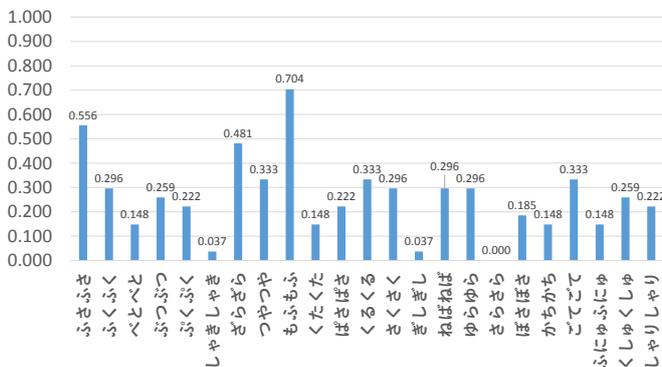


図 5: ユーザー評価の結果

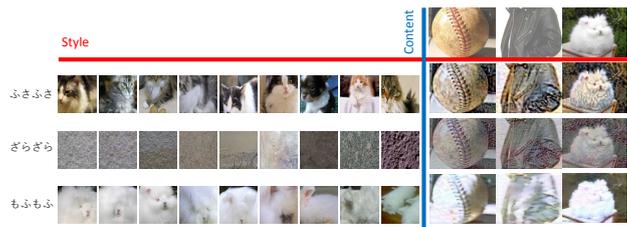


図 4: 正解率の高い質感変換結果例 (vote: 正解者の数)

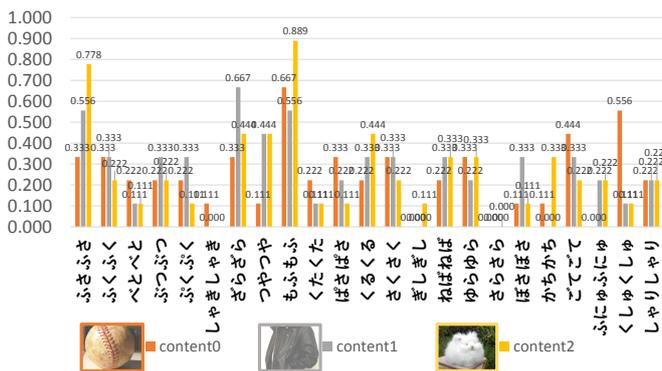


図 6: コンテンツ別ユーザー評価の結果

変換されており、自然な変換となっている。

したがって、コンテンツの持つ特徴を考慮したスタイルクラスタ構成により効果的なスタイル生成が可能になると思われる。

「ふさふさ」、「くしゅくしゅ」のスタイルクラスタでは全てのクラスタで似たスタイルのものが集まり、正解数にクラスタ毎の差はあまり見られなかった。それに対し、「ざらざら」、「もふもふ」のスタイルクラスタは3つのクラスタの内、1つのクラスタの正解率だけが非常に高くなっていた。「ざらざら」のスタイルクラスタは Cluster0 はパッチ状のスタイルとして有用性の高い画像があつまっており、それ以外のクラスタは類似した画像があまり集まっていなかったためだと思われる。「もふもふ」の Cluster0 は類似したスタイルが上位に集まっているが、アニメの画像など、スタイルとして好ましくない画像が多く、Cluster1 は「ふさふさ」などの他のオノマトペで収集された画像に似ていることから、選択されることが少なかったのではないと思われる。このクラスタの順序と正解率の違いからクラスタのスコア (式 6) は改良の余地があると考えられる。

また、これらの結果から、スタイル変換のクオリティについてはコンテンツ画像とスタイルの相性に強く依存していると考えられる。コンテンツが基調とするカラーによって背景と物体のどちらが大きく変化するかが変わっており、後者のコンテンツの正解率が高くなっている。例えば、図 4 の元となったスタイルクラスタでは、「ざらざら」は黒を、「もふもふ」は白を基調とした画像が集まっており、それぞれ同じ色を基調としたコンテンツを変換した画像 (「ざらざら」と content1, 「もふもふ」と content2) の正解率が高くなっている。

また、コンテンツとスタイルの間の共通構造 (しわ、毛並みなど) の有無によって生じるスタイル変換のクオリティの違いも影響していると思われる。例えば、同じく白を基調としたコンテンツであるにも関わらず、content0 と content2 の「ふさふさ」と「もふもふ」の正解率には大きな差がある。これらの単語では content0 では新たに毛の質感を追加しているのに対し、content2 では元より持っている質感を強調される形で



図 7: 「ふさふさ」による質感変換結果

## 5. おわりに

Neural Style Transfer アルゴリズムを用いた画像の任意質感生成を行うために、単語概念からスタイル表現を自動生成するシステムを構築し、複数の生成画像の評価をアンケートを用いて行った。その結果、視覚的概念への連想が容易であり、類似スタイル表現を多く持つクラスタが生成される単語では人間に評価されやすい画像が生成され、またそのクオリティはコンテンツのカラーやスタイルとの共通構造によって、大きく影響を受けることが分かった。

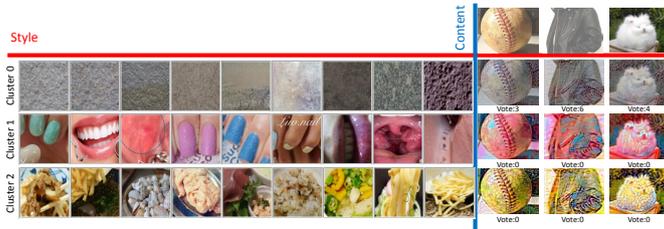


図 8: 「ざらざら」による質感変換結果



図 10: 「くしゅくしゅ」による質感変換結果



図 9: 「もふもふ」による質感変換結果

スタイル変換アルゴリズムの傾向から、変換に使用するスタイル画像はきめ細かな構造が多く、コンテンツ画像に類似する構造が有り、変換したい対象の領域のカラー傾向が共通していることが望ましいことが分かった。そこで、コンテンツのスタイル表現に依存したスタイルクラスターの構築や領域分割とスタイル画像の色変換を用いて対象領域のカラーをコンテンツに近づけることでより自然な質感変換が可能になると考えられる。

## 参考文献

- [1] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *Proc. of IEEE Computer Vision and Pattern Recognition*, 2016.
- [2] S. Matsuo and K. Yanai. Cnn-based style vector for style image retrieval,. In *Proc. of ACM International Conference on Multimedia Retrieval*, 2016.
- [3] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proc. of arXiv:1409.1556*, 2014.
- [4] C. Li and M. Wand. Combining markov random fields and convolutional neural networks for image synthesis. In *Proc. of arXiv:1601.04589*, 2015.
- [5] R. Novak and Y. Nikulin. Improving the neural algorithm of artistic style. In *Proc. of arXiv:1605.04603v1*, 2016.
- [6] G. Berger and R. Memisevic. Incorporating long-range consistency in CNN-based texture generation. In *Proc. of arXiv:1606.01286v1*, 2016.
- [7] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Proc. of arXiv:1603.08155*, 2016.
- [8] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *Proc. of arXiv:1603.03417v1*, 2016.
- [9] V. Dumoulin, J. Shlens, and M. Kudlur. A learned representation for artistic style. In *Proc. of arXiv:1610.07629v1*, 2016.
- [10] R. Datta, D. Joshi, J. Li, and J. Z. Wang. Studying aesthetics in photographic images using a computational approach. In *Proc. of European Conference on Computer Vision*, 2006.
- [11] L. Marchesotti and F. Perronnin. Learning beautiful (and ugly) attributes. In *Proc. of British Machine Vision Conference*, 2013.
- [12] Naila. Murray, De. Barcelona, L. Marchesotti, and F. Perronnin. Ava: A large-scale database for aesthetic visual analysis. In *Proc. of IEEE Computer Vision and Pattern Recognition*, 2012.
- [13] Keren. D. Painter identification using local features and naive bayes. In *Proc. of International Conference on Pattern Recognition*, 2012.
- [14] S. Karayev, M. Trentacoste, H. Han, A. Agarwala, A. Darrell, T. Hertzmann, and H. Winnemoeller. Recognizing image style. In *Proc. of British Machine Vision Conference*, 2013.
- [15] T. Y. Lin and S. Maji. Visualizing and understanding deep texture representations. In *Proc. of IEEE Computer Vision and Pattern Recognition*, 2016.
- [16] 幾田光, 大垣慶介, 小田桐優理. 畳み込みニューラルネットワークを用いたテクスチャ特徴量の混合に基づく自然なテクスチャ転写. 第 19 回 画像の認識・理解シンポジウム (MIRU 2016), 2016.