

深層学習による質感文字生成

成沢 淳史 下田 和 柳井 啓司
Atsushi Narusawa Wataru Shimoda Keiji Yanai

*1 電気通信大学 大学院情報理工学研究科
The University of Electro-Communications, Tokyo

In this work, we present a method to convert computer character fonts into natural fonts made of various kinds of materials such as ketchup, rope, and sands. To do that, we propose to combine a domain transfer method based on CycleGAN and a style transfer method. By the experiments, we confirmed that the proposed method was promising.

1. はじめに

文字にまつわる分野において深層学習のおかげで様々な研究タスクを考えられるようになってきている。情景文字認識のほか、古文書の解析や画像中のテキストを使ったイメージキャプションができるようになった。近年、コンピュータビジョン分野において画像生成タスクが盛り上がりを見せており、フォントの形状変換や新しいフォントの生成を実現するための応用が始まりつつある。フォント生成は中国語や日本語のように文字種が数千を超えるような言語でのフォント作成におけるコスト削減の点で有益な研究である。すでにフォントからフォントへの変換は文字の分野で重要なタスクとして様々な研究者が深層学習の手法を用いて試みている。本研究ではフォントからフォントへの変換タスクをさらに発展させ、身の回りに存在するパターンや共通のデザインが統一的に使われた画像からユニークな特徴をフォント画像に対して転写するタスクに挑戦する。こうした流れは、従来の文字を認識することが最終目標であった文字認識研究の枠を超えた新しい動きであり、文字にまつわる様々な研究ということで文字工学と呼称されている。

2. 実験目的

本論文では画像を扱う深層学習に基づいた文字画像の形状変換ないし生成によるフォント生成の最適な手法の探求を目的として行う。

日本語や中国語は数千におよぶ文字の種類があり、フォントの製作にはコストと時間が必要となる。したがって、画像生成技術により数種類のデザインパターンから日本語フォントを自動生成することが期待されている。しかしながら、深層学習を文字画像に利用する場合には生成結果の可読性の向上や少数サンプルからの学習といった課題が残っている。

深層学習による画像生成タスクの困難な点としては、深層学習固有の問題として多様かつサンプル数の多いデータセットが必要となる点が挙げられる。また、現状ではフォント間での変換には同じ対応する文字から変換パターンを学習するためデータセットの用意が難しく、フォント間以外への応用が困難である。そこで、本研究では以下の取り組みを行った。

1. Neural Style Transfer の導入による生成結果の可読性の改善
2. 質感パターンセットを作成し、フォントに対する質感の転写

3. 入力画像の工夫による少数サンプルからの学習

3. 関連研究

3.1 従来手法によるフォント画像生成

フォントにはベクトルフォントとビットマップフォントがある。一般的なベクトルフォント生成では文字を直線やカーブ、ハネなどのストロークの集合または部首のようなまとまったコンポーネントの組み合わせから成ると考え、文字をストロークにまで分解し、対応する変換先のストロークを割り当てる手法 [1, 2] やストロークやコンポーネントの組み合わせから漢字を作り出す手法 [3, 4] が提案されている。

そのため、現状では予めフォントのベクトル情報を利用してスケルトンを作成することでストロークへの分割が容易になり、漢字を構成する上で必要なストロークを特定、抽出しフォントの自動生成を行う。漢字を構成するために必要なストローク成分は研究されているが、漢字や記号毎にストロークの対応関係を記述するのは大変な労力となる。そこで深層学習の手法を用いることでフォント画像から自動でストロークを抽出し、変換対応をネットワークで学習させることで自動化を行うことができるようになる。

また、これまでの研究ではフォントの生成に絞った研究が多く見られたが、近年はフォントの自動生成よりもタイポグラフィを含めて飾り文字のようなアーティスティックな文字画像の生成 [5] に注目が集まる中で文字画像に対するテクスチャのような質感の付与のほか、文字以外のパターン画像に対する手法の適用が期待されている。本研究では深層学習の手法によりテクスチャのような質感の付与や文字画像からストローク間での対応関係を変換ネットワークで学習させることに挑戦する。

3.2 深層学習によるフォント生成

深層学習を用いた文字画像生成タスクの多くは中国圏の研究成果が多く著名な成果に Rewrite^{*1}, Zi2Zi^{*2} プロジェクトが存在する。Rewrite は Neural Style Transfer [6] をフォント画像の生成に適した変更を加えたプロジェクトである。Neural Style Transfer は二枚の画像を合成する手法であり、Rewrite の他にもフォントでの応用結果 [7] が報告されている。

さらに敵対学習の手法を利用した Zi2Zi がある。これは Pix2Pix [8] をベースに拡張を施したプロジェクトであり、画像

*1 <https://github.com/kaonashi-tyc/Rewrite>

*2 <https://kaonashi-tyc.github.io/2017/04/06/zi2zi.html>

から画像への変換を目的とした仕組みである。画像を特徴表現するためのエンコーダ、特徴表現から画像に復元するデコーダの変換ネットワークが使われており変換ネットワークの学習では変換元と変換先で同じ対応関係のペアを用意し学習を行う。

本研究ではこれらの研究をさらに発展させ、ペア画像無しでのクロスドメイン学習、フォント間以外の質感パターン画像セットからの学習に応用した。

4. 提案手法

画像セット名	サンプル数
ケチャップ文字	445
砂文字	483
紐文字	796

表 1: 質感画像セット概要



図 1: 質感パターン画像セット (左からケチャップ文字, 砂文字, 紐文字)

本研究ではフォント画像からケチャップ文字, 砂文字, 紐文字の3種類の質感パターン画像セット(図1)への変換実験を行う。本研究ではFast Style Transfer [9, 10]を参考に敵対学習のひとつであるクロスドメイン学習による手法 [11, 12, 13, 14]を組み合わせた, ネットワーク(図2)を考案した。クロスドメイン学習は変換ネットワーク部分 G, F(図3)とDiscriminator部分 D_x, D_y (図4)から成り, 順変換と逆変換を行うため, それぞれ2つずつ独立してネットワークが準備される。

変換ネットワークは入力画像の形状を保つため, 逆変換を行った画像とElement-Wiseでの比較によるCycle Loss (L_{cycle})に加え, Adversal Loss ($L_{adversal}$)がDiscriminatorによりフィードバックされ, 敵対学習を進めることで入力の形状に対してターゲットのデザイン, スタイルを転写することができる。

クロスドメイン学習の生成結果における可読性の改善のため, Style TransferのContent Loss ($L_{content}$)に注目し, また, 新しいデザインの獲得に繋がる結果が期待されるStyle Loss (L_{style})を導入している。Style Transferを導入するに辺り, 特徴抽出にはVGG 16のプレトレインモデルから行い, コンテンツレイヤーとスタイルレイヤーは図5に示す通り, コンテンツ特徴は1箇所, スタイル特徴は4箇所のレイヤーから抽出する。提案手法ではこれら4つのロスを統合しており, 変換ネットワークG, Fは式1のロスから学習される。それぞれの重み $\alpha, \beta, \gamma, \delta$ は実験的に決めており, 最終的な重みの設定は表2の通りとなっている。

$$L_{total} = \alpha L_{style} + \beta L_{content} + \gamma L_{adversal} + \delta L_{cycle} \quad (1)$$

Style Weight	3.00E+05
Content Weight	1.50E+00
Adversal Weight	4.50E+08
Cycle Weight	3.50E+12

表 2: ロスの重みの設定

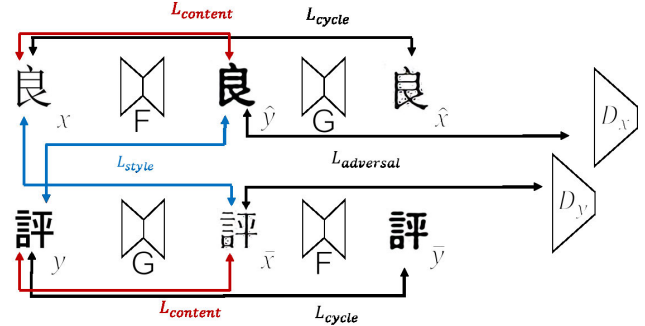


図 2: CycleGAN with Neural Style 概要

5. 実験

5.1 画像と実験環境

ケチャップ文字画像セットは文字の画像が多く, 砂文字には英字と簡単なひらがなに加えて絵の一部から構築している。また, 紐文字に至っては文字が含まれておらず, 紐で作られたアート画像からランダムで切り出しを行い手作業で整形を行った。さらに入力画像に関して, 画像中から十分なスタイル情報を得るために単一画像中に16文字を画像セットから選び, 配置を行い500枚として学習した。

5.2 生成例

Neural Style, Cycle GAN, 提案手法での生成例を図6に掲載する。フォントからそれぞれの質感画像セットへの変換が目的となるため, Cycle GANでは順方向への変換結果のみ掲載する。また, 表2とは異なり, Content Lossの重みが大きく設定されている。Cycle GANと提案手法において質感パターン画像への変換では文字の形状がはっきりしなかった生成画像がコンテンツロスの導入により改善される結果が得られた。特に紐文字のNeural Style (Style Loss + Content Loss)での生成はスタイルに重みを置いた場合に可読性を保持できないが, 提案手法では可能となっている。また, Style LossとContent Lossはトレードオフの関係にあるが, Adversal Loss, Cycle Lossが可読性や質感の転写を補う結果も得られた。

5.3 主観評価と客観評価

Style LossとContent Lossの導入の効果を確認するため, 様々なロスの組み合わせを検証した。図7ではContent Lossの有無により, 可読性が保持されていることが確認できるほか, Style + Cycle + Contentの組み合わせにAdversalを加えると生成画像のテクスチャが自然になることが視覚的にわかる。

また, 対象のパターン画像のスタイルをネットワークが獲得できているかをStyle Lossの平均(スタイル平均)により比較を行った。それぞれ画像セットから1枚のスタイル画像を選び, 生成した検証画像セットのStyle Lossを全て計り, 表3にまとめる。Neural Styleを使った場合のスタイル平均が最も小さい。本研究が目指すところは, この値とCycle GANをベースラインとして中間値をターゲットとしている。ケチャップ文字ではStyle Loss + Cycle Loss + Content Lossの組み合わせ

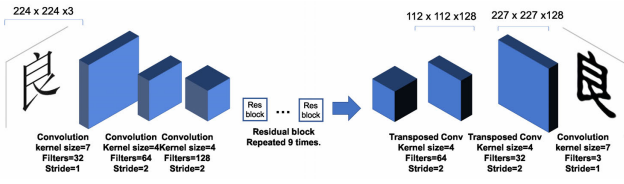


図 3: Encoder-Decoder Network 詳細

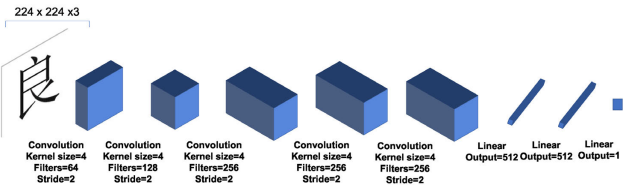


図 4: Discriminator 詳細

わせて達成しており, Adversal Loss を加えた場合には Cycle GAN より微減の結果となった。

一方, 砂文字では Neural Style と Cycle GAN でのスタイル平均に大差はなく, Style Loss を組み合わせたパターンは期待通りに作用しない一方で, Adversal Loss によるスタイル平均の減少の結果が得られた。最後に紐文字では, ケチャップ文字と同様の傾向にあるが, 複数のロスを組み合わせる場合うまくスタイルを獲得できる結果となった。

6. 考察

本提案手法では Adversarial Loss と Style Loss, Cycle Loss と Content Loss が競合する形になっている。しかしながら, 3つの画像セットを試したところ, それぞれでうまく機能する組み合わせとそうでないケースがあった。共通して Content Loss は入力形状の保持に働くが, Content + Style ではスタイル特徴をうまく合成できないケースがあった。この場合には, Cycle Loss を加えることで改善するケースがあったが, 考えられる要因として直線的なデザインの多い人工のフォント画像と輪郭線に歪みのあるパターン画像とで Gram matrix を通して共起関係を得られたかの差であると考えられる。Cycle Loss では2つの変換ネットワークを通すため, 直線部分に歪みが生じることで共起関係が得やすいことが予想される。Content Loss の重みを固定した状態で Cycle Loss を変化させるとスタイルの特徴が大きく混ざった視覚的に意味のある結果 (図 12) も実験より得られている。

また, Style Loss と Adversal Loss は性質が異なるものであると予想のもとから導入を行い実験を行った。実際には Adversal Loss が Style loss を改善するケースは砂文字画像セットでのみ結果が得られた。ケチャップ文字, 紐文字における提案手法ではスタイル平均が Neural Style 以上, Cycle GAN 以下となっている。Adversal Loss の重みを大きくすると背景のテクスチャの結果が改善する結果が砂文字より得られており,

	ケチャップ文字	砂文字	紐文字
Style + Content (Neural Style)	5.17E+06	2.77E+06	2.03E+07
Adversal + Cycle (Cycle GAN)	6.03E+06	3.11E+06	2.74E+07
Style + Cycle + Content	5.66E+06	3.36E+06	2.17E+07
Adversal + Style + Cycle	5.98E+06	2.69E+06	2.06E+07
Adversal + Style + Cycle + Content	6.00E+06	2.71E+06	1.99E+07

表 3: スタイル平均

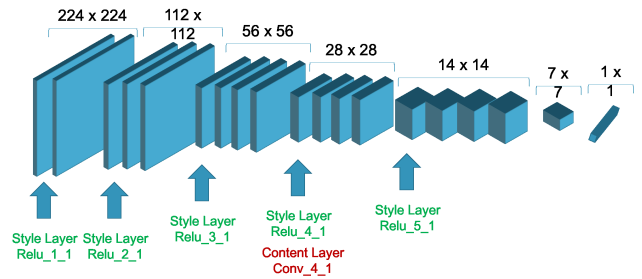


図 5: スタイルレイヤーとコンテンツレイヤーの選択



図 6: 生成結果 (Neural Style, Cycle GAN, 提案手法)

同等の効果は Content Loss により得られると思われる。Discriminator はテクスチャ, 形状などのディテールに対する制約なので, 自然な画像生成に繋がると考えられ, Content Loss は前景と背景に対する制約と考えられるが, 現状 Adversal Loss の方が自然な画像生成に繋がる印象を得ている。

様々なネットワーク構造で実験を行ったところ, Adversal Loss + Content Loss + Cycle Loss ないし, これに Adversarial Loss を加えた場合の組み合わせが最も良く, 特に Cycle GAN に Content Loss を導入した場合には可読性が高い可能性で改善すると考えている。特に図 8 に示すように, ひらがなやカタカナなどの画数の少ない簡単な文字ほど, 可読性の高い変換が実現できている。

7. 今後の課題

生成結果の評価を行うためには質感 (デザイン) と可読性の二つの指標で測る必要がある。実験では質感の転写がなされているのかを Style Loss を使って比較を行った。今後, この Style Loss による問題点があるのかを議論したいと考えている。また, 文字の可読性については今回測ることができなかった。有力な手段として万能な文字分類器を作成し, 文字の分類率により比較を行うのが最も良いと思われる。

今回の手法はパッチレベルでの変換対応がうまくいかない場合への対処として裏付けがなされていない。文字は弾性を持っているため, 形状が大きく変わった場合でも可読性が変わらない

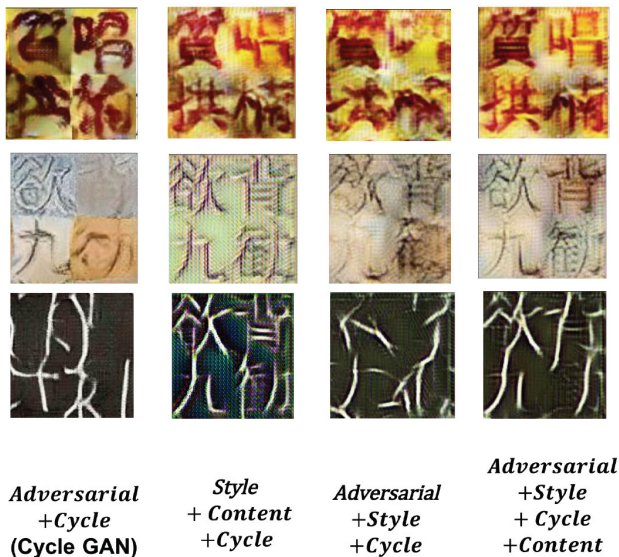


図 7: ロスの組み合わせの比較



図 8: シンプルなひらがな・カタカナと、複雑な漢字の変換例。

い場合がある。スタイルに適した形状まで予め入力画像を摂動させることができれば、変換対応の学習に繋がり、ストロークの消失に対するアプローチになる可能性がある。

謝辞: 本研究に対しご助言、ご議論をして下さった九州大学内田誠一先生に感謝いたします。また本研究は JSPS 科研費 17H06100 基盤 (S)「機械可読時代における文字科学の創成と応用展開」の助成を受けたものです。

参考文献

- [1] A. Zong and Y. Zhu. Strokebank: Automating personalized chinese handwriting generation. In *Advancement of Artificial Intelligence*, pages 3024–3030, 2014.
- [2] T. Miyazaki, T. Tsuchiya, Y. Sugaya, S. Omachi, M. Iwamura, S. Uchida, and K. Kise. Automatic generation of typographic font from a small font subset. In *arXiv preprint arXiv:1701.05703*, 2017.
- [3] J. Lin, C. Hong, R. Chang, Y. Wang, S. Lin, and J. Ho. Complete font generation of chinese characters in personal handwriting style. In *Computing and Communications Conference (IPCCC), 2015 IEEE 34th International Performance*, pages 1–5. IEEE, 2015.
- [4] X. Songhua, L. Francis C.M., C. Kwok-Wai, and P. Yunhe. Automatic generation of artistic chinese calligraphy. In *IEEE Intelligent Systems*, 2005.
- [5] Y. Shuai, L. Jiaying, L. Zhouhui, and G. Zongming. Awesome typography: Statistics-based text effects

transfer. In *Proc. of IEEE Computer Vision and Pattern Recognition*, 2016.

- [6] 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.
- [7] A. Gantugs, B. K. Iwana, A. Narusawa, K. Yanai, and S. Uchida. Neural font style transfer. In *International Conference on Document Analysis and Recognition*, 2017.
- [8] I. Phillip, Z. Jun-Yan, Z. Tinghui, and E. Alexei. Image-to-image translation with conditional adversarial networks. In *Proc. of IEEE Computer Vision and Pattern Recognition*, 2017.
- [9] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Proc. of European Conference on Computer Vision*, 2016.
- [10] A. Vedaldi D. Ulyanov, V. Lebedev and V. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *arXiv:1603.03417v1*, 2016.
- [11] Z. Jun-Yan, P. Taesung, I. Phillip, and E. Alexei. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proc. of IEEE International Conference on Computer Vision*, 2017.
- [12] Z. Yi, H. Zhang, P. Tan, and M. Gong. Dualgan: Unsupervised dual learning for image-to-image translation. In *Proc. of IEEE International Conference on Computer Vision*, pages 2849–2857, 2017.
- [13] T. Kim, M. Cha, M. Kim, J. K. Lee, and J. Kim. Learning to discover cross-domain relations with generative adversarial networks. In *International Conference on Machine Learning*, 2017.
- [14] Y. Taigman, A. Polyak, and L. Wolf. Unsupervised cross-domain image generation. In *International Conference on Learning Representation*, 2017.