

敵対的生成ネットワーク (GAN) を用いた似顔絵生成手法の検討

中島悠輔¹ 坂内祐一²

概要: 似顔絵は人物の外見・特徴をとらえて、デフォルメして描いた人物画である。現在、顔画像を似顔絵に変換する研究では、非教師学習を用いた変換手法や、それぞれパーツ毎に変換を行う手法の研究が行われている。しかし、プロのイラストレーターの個性を反映させるような研究は発表されていない。そこで、本研究ではプロのイラストレーターの個性を反映する為に、どのネットワークを用いれば良いか検討を行う。具体的には、pix2pix, CycleGAN ペア, CycleGAN 非ペア, Cyclepix の4つのネットワークを検討していく。違いとしては主に損失関数となっている。pix2pixは訓練データと生成データの誤差を取るが、CycleGANでは入力データと生成データを更に変換させた再変換データの誤差を取ることが主な違いである。Cyclepixは両方の誤差を取っている。また、pix2pix, CycleGAN ペアはDiscriminatorの入力が生成データのみとなっていることも違いである。実験結果として、CycleGAN 非ペアと Cyclepix の評価が高いことが分かった。このことからDiscriminatorの入力が生成データのみ、Cycle Consistency Loss を利用することで精度が高い似顔絵を生成することに有用である。

キーワード: 似顔絵, 敵対的生成ネットワーク, Deep Learning, GANs

A generation method of cartoon portrait using Generative Adversarial Networks

YUSUKE NAKASHIMA^{†1} YUICHI BANNAI^{†2}

1. はじめに

似顔絵は様々な場面で使われている。SNS, ブログゲーム, アプリのプロフィールやアバターといった用途でプライバシーの観点から顔写真ではなく、自分のイメージを表現する為に似顔絵を用いることが良く見られる。

しかし、似顔絵を描くにはセンスとスキルを必要とされており、これらを持って無い人は似顔絵を描くことが出来ない。従って自動的に似顔絵生成を行うシステムを作成することの重要性が高まっている。

従来の研究では、写真加工.com[1]というWebサイトの色や明るさを減らし、線を加えるエッジのポストリゼーションと呼ばれる手法で画像処理による似顔絵の生成システムや呉ら[2]により研究されている顔画像から特徴点を取得しニューラルネットを用いて顔の個性に対応する非線形変換を施し、線画ベースの似顔絵を自動的に生成するシステム等が提案されている。しかし、写真加工.comの場合はパラメータを手動で調整しなければならないこと、写真を似顔絵風に変換する為トレース画となること、写真自体の精度が悪いと似顔絵生成の精度が悪くなる場合がある。呉らの研究では、ニューラルネットによって得られるパラメータから似顔絵生成していることからパターンとなってしまう同じパターンが出来上がってしまう問題がある。

本研究では、Deep Learningにおける画像生成に使われるGAN (Generative Adversarial Network, 敵対的生成ネットワーク) を用いて、プロのイラストレーターが描いたような似顔絵を少ない訓練データで生成する手法を検討していく。顔画像とその顔画像を元にプロのイラストレーターが描いた似顔絵のペア用意し学習させた後、顔画像から似顔絵の生成を行う。顔画像と似顔絵のペアは男女年代がそれぞれ分布するように訓練データを作成した。GANのネットワークであるpix2pix, CycleGANの非ペアとペア、本研究で提案する手法であるCycleGANとpix2pixの手法を足し合わせたネットワーク(以下Cycle pix)を利用して似顔絵画像生成を行い比較する。似顔絵の訓練データが少ないことを想定し、90枚と189枚の場合について画像生成結果を用いて評価実験を行い、顔画像から似顔絵の最適な生成手法を検討する。

2. 関連研究

2.1 Generative Adversarial Nets (GAN)

I.J.Goodfellowら[3]は、Deep Learningにおける効率的に生成モデルを訓練するための手法であるGANを提案した。GANではGenerator(生成モデル)とDiscriminator(識別モデル)を作成し訓練を行い、Generatorは訓練データに似ているデータを生成し、Discriminatorは訓練データである確

¹ 神奈川工科大学大学院 情報工学専攻
Kanagawa Institute of Technology
² 神奈川工科大学 情報メディア学科
Kanagawa Institute of Technology

率を出力する識別器となる。Generator は生成したデータを Discriminator が訓練データと識別するように、Discriminator は訓練データに対して訓練データと識別し、Generator が生成したデータに対して Generator が生成したデータと識別するように訓練を行う。このように Generator と Discriminator を敵対的に訓練することで、互いに競い合っ て精度を向上させていく手法となる。問題としては学習が安定しないことが示されている。

2.2 Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Network (DCGAN)

A.Radford ら[4]が、DCGAN (Deep Convolutional Generative Adversarial Network) の使用を提案した論文であり、教師なし学習による画像生成の精度の向上を目的としている。要点として畳み込みにおける pooling の排除、Batch Normalization による Generator 学習の効率化、全結合層の排除、Generator の出力層に tanh 関数を採用し、Discriminator の全ての活性化関数を LeakyReLU に変更を行っている。結果として GAN の学習を安定させることに成功している。本研究では Generator、Discriminator のネットワーク構造に DCGAN が提案している手法を用いている。

2.3 Image-to-Image Translation with Conditional Adversarial Networks (pix2pix)

P.Isora ら[5]は GAN を利用した画像生成アルゴリズムの一種で、二つのペアの画像から画像のドメインの関係を学習することで、1枚の画像からその関係を考慮した補間をしてペアの画像を生成する技術である pix2pix を提案した。pix2pix では、Conditional GAN (条件付き GAN) を採用している。基本方針は DCGAN を用いており、Generator として U-NET を採用している。また、Discriminator には PatchGAN を用いており、全体を見るのではなく局所に注目し訓練データか Generator が生成したデータかを見分けている。しかしながら、訓練データ量が大きい必要があることや教師画像をペアで用意しなければならないことが示されている。

2.4 Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks (CycleGAN)

J.Zhu ら[6]は 2.3 の pix2pix と同様、GAN を利用した画像生成アルゴリズムの一種である CycleGAN を提案した。画像のドメインを定めてドメインごとに画像を集めて訓練データとする。対になる訓練データを用意する必要がないのが特徴となる。ドメインごとの画像の集合を X,Y とし、それらに対して X から Y、Y から X の変換を行う Generator を用意する。加えて、双方に対応する Discriminator も二つ用意する。提案手法では GAN に用いられている Adversarial Loss とこの論文で提案された Cycle Consistency Loss の誤差 (loss) を学習させることにより非ペアの画像変換を可能としている。しかし、非教師学習の為、pix2pix に比べ精度が悪いことが示されている。

2.5 APDrawingGAN: Generating Artistic Portrait Drawings from Face Photos with Hierarchical GANs

R.Yi ら [7] は似顔絵を描く為の GAN である APDrawingGAN を提案した。画像全体の変換を行う Global net と右目、左目、鼻、口、髪の変換を行う Local net に分かれており、最後に Fusion net で結合を行うアーキテクチャを採用している。また、似顔絵によくみられる小さなずれを許容し、新しく Line-promoting distance transform loss を提案した。それらにより、少ない訓練画像で高精度な似顔絵の生成に成功した。

3. 実験を行ったネットワーク

関連研究で挙げた教師有りで変換を行える pix2pix、教師無しで変換を行える手法である CycleGAN 非ペアを採用した。pix2pix は航空地図からマップや線画ベースの鞆からカラー画像の鞆の変換に成功したことが報告されており、教師有の場合の基本的な変換手法であることから採用した。CycleGAN 非ペアは教師無しで過学習を起しにくく、高精度な変換が行えることから採用した。しかし pix2pix に比べて CycleGAN 非ペアの方が過学習を起しにくい、教師無し学習の為訓練データ通りの生成が出来ないことがある。その為、pix2pix の Discriminator の入力をペアにすることで教師付けた CycleGAN ペアに加えて、pix2pix のピクセル誤差を加えることで教師付けた Cyclepix を提案することとした。本章では、実験を行ったネットワークである pix2pix、CycleGAN ペア、CycleGAN 非ペア、提案する手法の Cyclepix について述べる。

3.1 ネットワーク

GAN は画像を生成する Generator と訓練データか Generator が生成したデータ化を識別する Discriminator から構成される。

3.1.1 pix2pix

pix2pix は教師あり学習でドメイン変換を行うネットワークである。訓練データに二つのペア画像を用意し、そのペア画像から画像間の関係を学習することで、1枚の画像からその関係を考慮した補完をしてペアの画像を生成するネットワークである。pix2pix の概要図を図 1 に示す。Generator は顔画像から似顔絵を生成する。Discriminator は顔画像と Generator が生成した似顔絵画像 or 訓練データの似顔絵画像をペアとして入力し、Generator が生成した画像か訓練データかを識別する。

3.1.2 CycleGAN ペア

CycleGAN ペアは CycleGAN が提案された論文で発表された内容を元に Discriminator の入力をペアに変更したネットワークである。Generator と Discriminator を一つずつ追加し、ドメイン変換が出来る手法で pix2pix のように 1対1の写像関係ではなく、どちらからも変換が出来る手法である。図 2 に CycleGAN ペアの概要図を示す。pix2pix との違いと

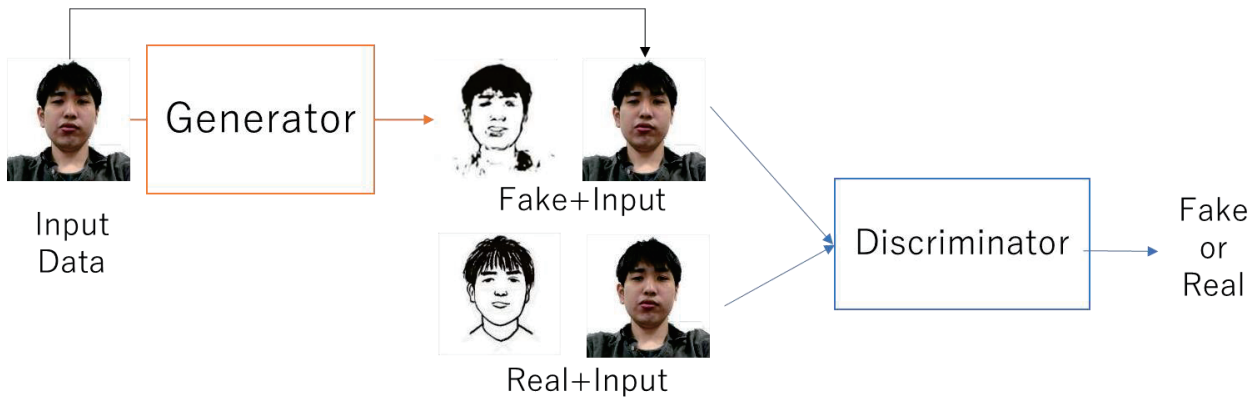


図 1 pix2pix の概要図. Generator は入力データである顔画像から似顔絵の生成を行う. Discriminator は似顔絵と入力データである顔画像ペアにして入力し, Generator が生成したデータか, 訓練データかを識別し, 識別結果を出力する.

して顔画像から似顔絵に変換する Generator1 に加えて似顔絵を顔画像に変換する Generator2 を追加している為, 再変換を行うことである. また, Generator2 の学習を行う為に生成した顔画像と Generator2 の入力である似顔絵をペアにして識別させる Discriminator2 が追加されている.

3.1.3 CycleGAN 非ペア

画像のドメイン変換を教師なし学習で行うことができる画像生成のネットワークである. ペアで学習しない為, ペアに比べて対訳が無くても画像と画像の変換が出来る手法である. これまで述べてきた手法より過学習を起しにくく, 少ない訓練データの場合でも高い精度の画像生成が望めるネットワークである. 図 3 に CycleGAN 非ペアの概要図を示す. 制約が緩く過学習を起しにくいメリットがあるが, ペアの条件付けがないために同じ人の似顔絵ではなく他人の似顔絵を生成するよう学習してしまう場合があ

り, メガネが無い人にメガネを描いてしまう場合や, ヒゲが無い人にヒゲを描いてしまう場合が生じるデメリットがある.

3.1.4 Cyclepix

今回提案するネットワークである. CycleGAN 非ペアを元に改良したネットワークであり, CycleGAN 非ペアで起こっている訓練データ通りにうまく学習できない問題を解決する為に提案した. CycleGAN ペアのように Discriminator で教師付けるのではなく, 生成した似顔絵と訓練データの似顔絵のピクセル誤差を加え, 教師づけたモデルである. ネットワークの概要図は CycleGAN の非ペアと変更が無い.

3.2 ネットワーク構造

本研究では, CycleGAN のネットワーク構造の構成を用いており, Generator と Discriminator を全てのネットワークにおいて同じ構造を用いている.

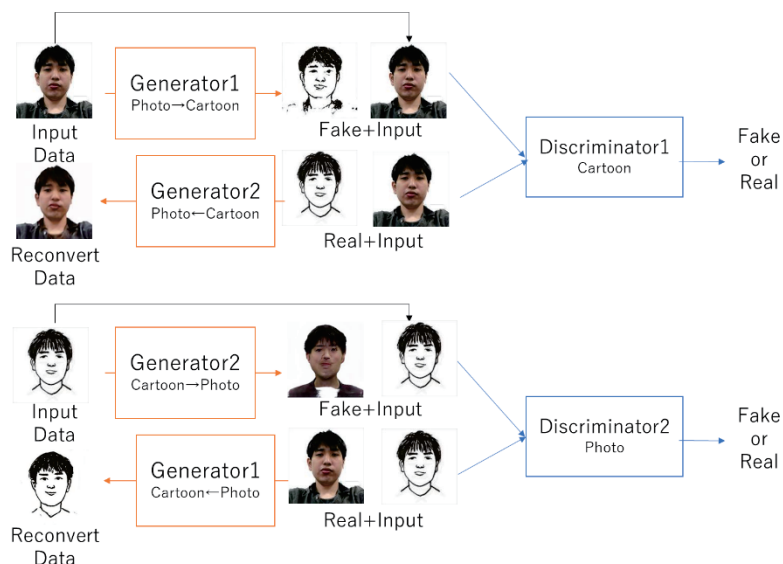


図 2 CycleGAN ペアの概要図. 顔画像から似顔絵を生成し, 顔画像に再変換を行う学習と似顔絵から顔画像を生成し, 似顔絵に再変換を行う学習がある. 生成した似顔絵から顔画像の生成を行う Generator2 を追加し, 再変換出来るようにしたネットワーク構造である. また, Generator2 を学習させる為に生成した顔画像と入力の似顔絵をペアにして Discriminator2 に識別させ, 結果を出力する.

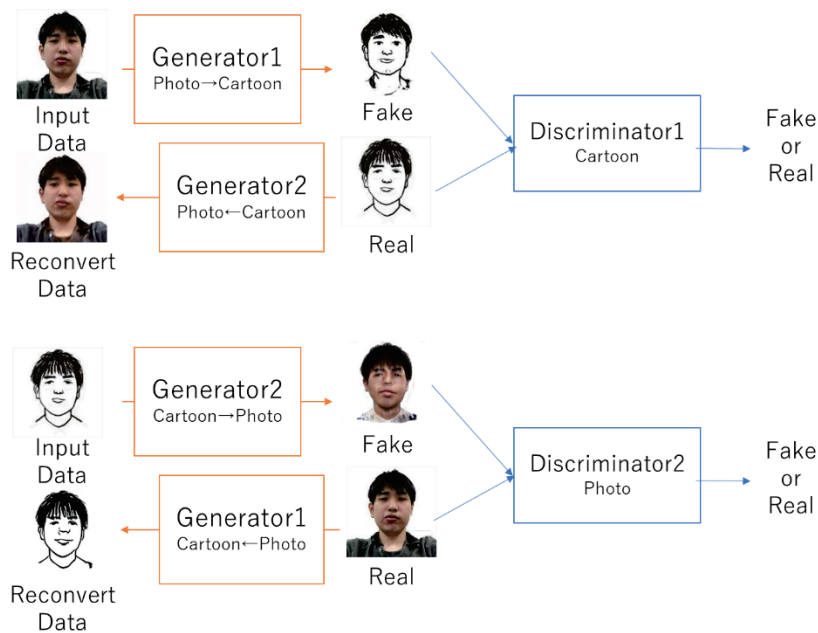


図 3 CycleGAN 非ペア, Cyclepix の概要図. 顔画像→似顔絵の学習を行う場合の概要図を示している. Discriminator の入力に Generator が生成した似顔絵のみとなっている.

3.2.1 Generator

Generator の構造図を図 4 に示す. Generator は顔画像の特徴を抽出する Encoder, 顔の特徴から似顔絵の特徴へと変換を行う Translation, 似顔絵の特徴から画像の形へと復元を行う Decoder の構造をしている. Encoder 部分では畳み込み層を 3 回用いて特徴マップを生成している. スライドによる畳み込みを用いてサイズを圧縮している. Translation では, ResNet[8]と言われる Residual Block を 9 個用いて変換を行う. Residual Block には Skip connection があり, 畳み込みを 2 回行った結果と入力を Skip connection を用いて足し合わせたものを出力する手法で, 勾配消失問題に強いと言われている. Decoder では 2 回転置畳み込み層を用いて画像のサイズまで復元させた後, 畳み込みを行って似顔絵画像を生成する. 活性化関数には出力層以外は Relu, 出力層のみ tanh を用いている. Batch Normalization は Instance Normalization を用いている.

3.2.2 Discriminator

Discriminator は, 入力に生成した画像のみを取る場合と生成した画像と Generator の入力の顔画像のペアを取る場合がある為, 入力のチャンネル数が変わることがある. Discriminator の構造は 2.4 節で述べたとおり PatchGAN といわれる手法を採用している. 3 回ストライドによる畳み込みを用いてサイズを圧縮しながら, 1 次元の畳み込みを行って, 32x32x1 の真偽値を出力させている.

3.3 損失関数

学習を行う上での手法やハイパーパラメータは CycleGAN を参考に行っている. 損失関数はネットワーク毎に異なっている. Generator の損失関数は各ネットワークで異なる. まず pix2pix では Generator が生成したデータと訓練データの似顔絵のピクセル誤差を取る L1 Loss, Discriminator に偽物と見抜かれたことを示す Adversarial Loss の重み付き和である. Adversarial Loss の式 1 は以下

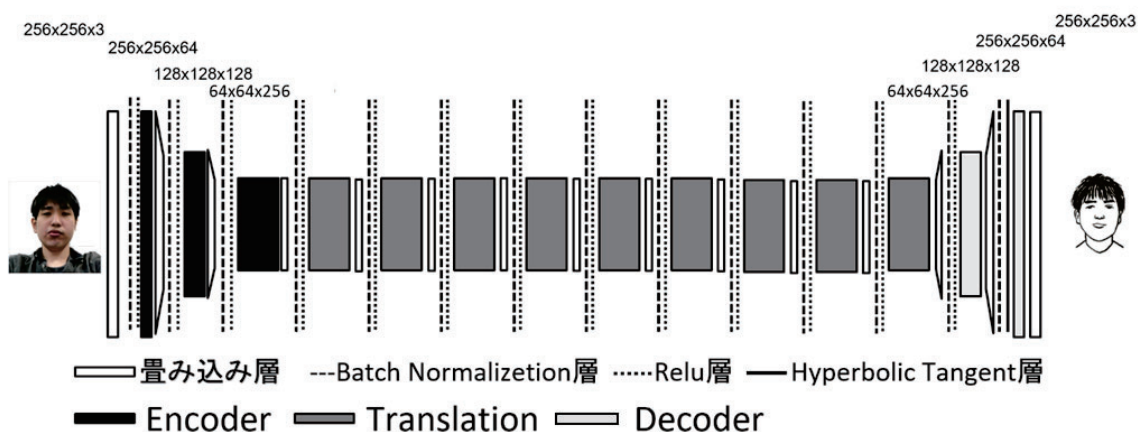


図 4 Generator のネットワーク構造

になっている。

$$L_{Adv} = -\mathbb{E}_{x \sim p_x(x)}[\log D(G(x))] \quad (1)$$

G, D はそれぞれ Generator, Discriminator の関数, x は訓練データ. 実験では, $\lambda_{Adv}=1.0$, $\lambda_{L1}=10.0$ とした.

$$L_G = \lambda_{Adv} * L_{Adv} + \lambda_{L1} * L_{L1} \quad (2)$$

次に CycleGAN である. こちらはペア, 非ペア共に同じ損失関数を用いている. 生成したデータを更に再変換させたデータと入力データの顔画像との誤差を取る Cycle Consistency Loss, Adversarial Loss の重み付き和である. 実験では $\lambda_{Adv} = 1.0$, $\lambda_{Cycle} = 10.0$ とした.

$$L_G = \lambda_{Adv} * L_{Adv} + \lambda_{Cycle} * L_{Cycle} \quad (3)$$

次に Cyclepix である. L1 Loss, Adversarial Loss, Cycle Consistency Loss の重み付き和である. 実験では $\lambda_{Adv} = 1.0$, $\lambda_{L1} = 2.5$, $\lambda_{Cycle} = 10.0$ とした.

$$L_G = \lambda_{Adv} * L_{Adv} + \lambda_{L1} * L_{L1} + \lambda_{Cycle} * L_{Cycle} \quad (4)$$

Discriminator の損失関数は全てのネットワークにおいて以下の式で示される Adversarial Loss のみである.

$$L_{Adv} = \mathbb{E}_{x \sim p_{data(x)}}[\log D(x)] + \mathbb{E}_{z \sim p(z)}[\log(1 - D(G(z)))] \quad (5)$$

3.4 最適化手法と学習

Optimizer には Adam[9] を使用, $\beta_1 = 0.5$, $\beta_2 = 0.999$ とした. 学習率は Generator, Discriminator 共に 0.0002 とした. 各 iteration につき Generator, Discriminator 共に 1 度ずつ更新し, ネットワーク毎にある程度精度良く生成できたと思われる段階で学習を打ち切った.

4. 実験

ネットワークの違いを評価するため, 顔画像と似顔絵画像を用いてモデルの学習を行った. また, 評価を行うため被験者を募り評価実験を行った.

4.1 訓練データ

顔画像とプロのイラストレーターが描いた似顔絵のペアを訓練データとする. 最初の段階として老若男女偏らないように女性の 0-40 歳, 40 歳以上, 男性の 0-40 歳, 40 歳

以上に分けてそれぞれ 25 枚ずつになるように集め, 合計 100 枚とし訓練データ 90 枚, テストデータ 10 枚とした. テストデータの選択基準として各区分から均等になるようにした. 平均すると 2.5 枚になる為, 女性の 0-40 歳, 男性の 40 歳以上から 2 枚と女性の 40 歳以上, 男性の 0-40 歳から 3 枚とした. その後訓練データに足りない髭, 眼鏡, 白髪等の特徴がある画像を 99 枚追加して, 189 枚の訓練データを準備した. 実験では 90 枚で訓練した場合, 189 枚で訓練した場合の 2 つを条件としている. 拡大縮小, 左右反転の画像拡張処理を行った. また画像は 300x300pixel で読み込み, 256x256pixel でトリミングを行っている.

4.2 生成結果

各ネットワークが生成したテストデータの結果を図 5 に示す. また基準となる手法として AutoEncoder[10]の結果を示す. 手法ごとに上段が 90 枚, 下段が 189 枚で学習させたモデルを用いている. AutoEncoder, pix2pix では似顔絵が顔のパーツの構成に失敗しており, CycleGAN ペアでは良い似顔絵もあるが失敗しているケースも多い. CycleGAN 非ペア, Cyclepix では比較的精度が高い似顔絵が生成出来ているが, CycleGAN 非ペアでは 189 枚より 90 枚の方が良い場合もあることが分かった. ネットワークの違いを比較する上で学習が訓練データ通りに上手くいっているかを確認した. AutoEncoder, pix2pix, CycleGAN ペアはほぼイラストレーター通りの出力が出来ているが, CycleGAN 非ペアはイラストレーター通りになっていないことが分かった. CycleGAN 非ペアは眼鏡を付けていないはずなのに眼鏡を描いていたり, 目や鼻といった各特徴の描き方が違っていたり髭が無いのに髭を描いていることが分かった. それに対し Cyclepix ではほぼイラストレーター通りの特徴を再現出来ていることが分かった. このことから AutoEncoder, pix2pix, CycleGAN ペアは訓練データに適合し過ぎていて, テストデータの精度が悪い過学習を起こしているのに対し CycleGAN 非ペア, Cyclepix は訓練データに適合し過ぎず, テストデータも一定の精度で生成出来ていることと推測さ

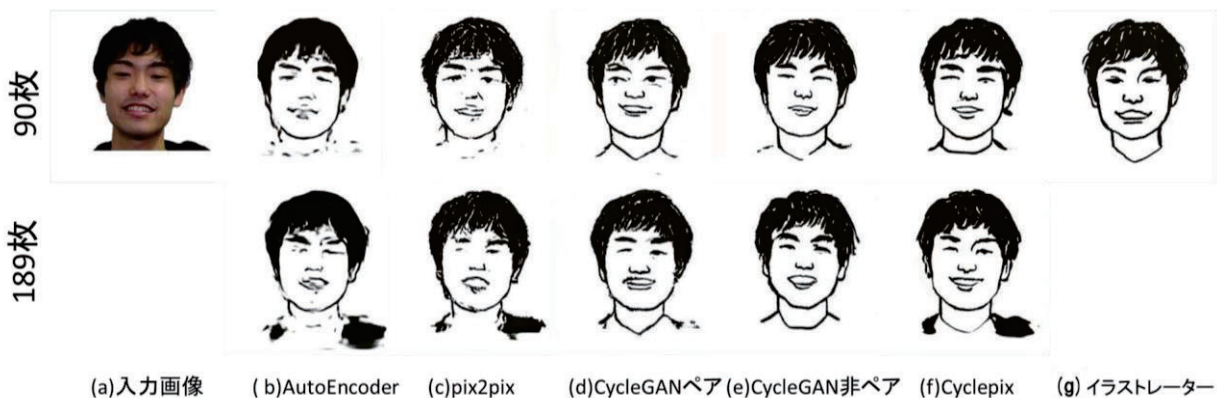


図 5 テストデータの生成結果. それぞれ上段が 90 枚, 下段が 189 枚で学習させたモデルを用いている. AutoEncoder, pix2pix では顔のパーツの構成に失敗しており, CycleGAN ペアでは良い似顔絵もあるが失敗しているケースも多い. CycleGAN 非ペア, Cyclepix は比較的良い精度で似顔絵が生成出来ている.

れる。

4.3 類似度の評価実験

GANの生成された似顔絵について、誤差の値をそのまま評価に用いることは適切ではない、その為イラストレーターが描いた似顔絵と各ネットワークが生成した似顔絵の類似度を評価として、被験者による評価実験を行った。評価実験に用いたデータは尺度データに用いたプロのイラストレーターが描いた似顔絵と各ネットワークが生成した似顔絵であるイラストレーターが描いた似顔絵と8つの生成結果を表示し各ネットワークが生成した似顔絵がプロのイラストレーターが描いた似顔絵と類似しているか、似ていない「-2」、やや似ていない「-1」、どちらでもない「0」、やや似ている「1」、似ている「2」とした5点法と順位法を用いて評価させた。大学3、4年生の計10名に対し、上記の評価実験を行った。表示する選択肢はランダムに場所を変更している。

4.4 評価実験の結果

結果を表1に示す5点法と順位法の結果からCycleGAN非ペア189枚が1.22、Cyclepix90枚が1.04、Cyclepix189枚が0.97、CycleGAN非ペアが0.4であることから似ていると判断されたことが分かる。CycleGANペアとpix2pixの値はマイナスの値となり、似ていないと評価された。また順位付けの結果からもCycleGAN非ペアとCyclepixが生成結果として優れており、1番似ていると判断されたものがCycleGAN非ペアであることが分かった。

表1 評価実験の結果

ネットワーク名	5点法	順位法
pix2pix 90枚	-1.65	7.15
pix2pix 189枚	-1.42	6.72
CycleGAN ペア 90枚	-0.85	5.63
CycleGAN ペア 189枚	-0.66	5.12
CycleGAN 非ペア 90枚	0.40	3.59
CycleGAN 非ペア 189枚	1.22	2.18
Cyclepix 90枚	1.04	2.68
Cyclepix 189枚	0.97	2.91

5. 考察

CycleGAN非ペアとCyclepixの評価が高かった理由として、顔の各パーツがぼやけておらずはっきりとした高精度な似顔絵が生成出来ていたからだと考えている。原因として、pix2pixやCycleGANペアは過学習が起きていたと考えられる。Discriminatorの入力がペアであることから、Adversarial Lossに条件付けがされていたことが主な理由としてあげられる。また、CycleGANとCyclepixに用いられているCycle Consistency Lossについてpix2pixに比べてCycleGANペアの精度が良くなっていることからL1 LossよりもCycle Consistency Lossが生成精度を高めることに寄

与していたと考えられる。pix2pixの過学習が起きていたことは似顔絵のようなデフォルメを行う場合に、輪郭や各特徴の形が変わる場合には入力と出力の整合性が取れず精度が低くなったためと考えられる。CycleGAN非ペアの189枚が上手くいった理由としては、訓練データの生成結果において、本来眼鏡を付けていない人に眼鏡、ヒゲが無い人にヒゲを描いたなどの例が見られたことから、教師無し学習の過学習を起こしていなかったことが考えられる。このことからデータを増やしていくとピクセル誤差で教師付けを行えるCyclepixの方が各特徴のイラストレーターの描き方の個性を反映できると考えられる。以上をまとめると現状としては訓練データが少ない枚数であることからヒゲ、フレームレスの眼鏡といった特徴を反映出来ない場合やGeneratorに対してDiscriminatorが強すぎてしまうといった問題があり、訓練データの枚数を増やすことやヒゲ、フレームレスの眼鏡といった特徴を抽出することが出来るようにネットワーク構造や画像サイズを変更すること。Discriminatorの学習係数を低く設定するか、Discriminatorの学習回数を減らすことによって弱める必要があると考えている。

参考文献

- [1] “似顔絵作成ソフト | 写真加工.com”, <http://www.photo-kako.com/likeness.cgi>, (参照日 2019-12-16)
- [2] 呉玉珍, 榎本誠, 川村春美, 大谷淳: 顔画像からの線画ベースの似顔絵自動生成システムにおける主観的識別に関する検討, FIT2014 第13回情報科学技術フォーラム, 第3分冊, pp.247-248, (2015)
- [3] I. Goodfellow. NIPS 2016 tutorial: Generative adversarial networks. arXiv preprint arXiv:1701.00160, 2016.
- [4] A. Radford, L. Metz, and S. Chintala.: Unsupervised representation learning with deep convolutional generative adversarial networks. In ICLR, 2016. 2
- [5] Isola, Phillip, et al. "Image-to-image translation with conditional adversarial networks." arXiv preprint arXiv:1611.07004 (2016).
- [6] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros.: Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, in IEEE International Conference on Computer Vision (ICCV), (2017)
- [7] Ran Yi, Yong-Jin Liu, APDrawingGAN: Generating Artistic Portrait Drawings from Face Photos with Hierarchical GANs, CVPR '19, pages10743-10752, 2019.
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In CVPR, 2016
- [9] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- [10] G. E. Hinton and R. R. Salakhutdinov.: Reducing the dimensionality of data with neural networks. Science, 313(5786):504–507, 2006.