# 輪郭線と色からなるカラースケッチからのフォトリアルな画像生成

落合 晃汰<sup>†</sup> 青野 雅樹<sup>††</sup>

† 豊橋技術科学大学 博士前期課程情報・知能工学専攻 〒 441-8580 愛知県豊橋市天伯町字雲雀ヶ丘 1-1
†† 豊橋技術科学大学 情報・知能工学系 〒 441-8580 愛知県豊橋市天伯町字雲雀ヶ丘 1-1
E-mail: †{ochiai.kota.yx,masaki.aono.ss}@tut.jp

**あらまし** 本研究では, pix2pixHD をベースとした手法を用いて, 輪郭線と色からなるカラースケッチからフォトリア ルな画像を生成する手法を提案する. 近年の深層学習を用いた画像生成の進歩はセグメンテーションマップやテキス トからの高品質な画像生成を可能にしたが, これらの手法ではラベルや言語に依存するために汎化性能やユーザビリ ティが低く, 広く実用化されるまでには至っていない. そこで本研究ではラベルや言語に依存しない入力形式として 輪郭線と色からなるカラースケッチ画像を用いてフォトリアルな画像を生成する手法を提案する. さらに生成モデル として Diffusion Model に基づく U-Net と DPT-Large を導入した手法を提案し, FID 評価尺度と KID 評価尺度, アン ケートによる評価にて精度改善したことを報告する.

キーワード 深層学習, GAN, 画像生成

# 1 はじめに

近年,深層学習による画像生成の活発化により非常に高品質 な画像が生成できるようになりつつある. 条件を設けない単純 な生成だけでなく、条件を与えた画像生成も非常に高品質に行 えるようになっている. 条件をユーザが与えることで欲しい画 像を生成しやすくなるため, 画像素材の作成や画像制作の補助 ツールとしての応用が可能となる.また、高度なスキルを用い ることなく画像生成が可能となるほか, 従来の画像制作ソフト と比較して短時間での画像制作の実現が考えられる.画像生成 時に設ける条件としては、セグメンテーションマップやテキス トを条件とした画像生成に関する研究が活発に行われている. セグメンテーションマップは画像のピクセル一つ一つに対して, 何が写っているかと言った、 ラベルやカテゴリを対応付けた画 像である.この方法は、ラベルやカテゴリ情報を画像に指定す るだけで良いため非常に簡単に画像生成が可能となるという利 点がある.一方で欠点も存在し、ラベルとして事前に設定して 学習させた物しか生成することができないという問題がある. さらに、 ラベルをどのような基準でどこまで分類すればいいか という点も問題となる. 例えば, 猫の画像を生成したいときに 「動物」というラベルしか存在しない場合, 猫の画像が生成され るとは限らない.しかし、数多存在するすべての動物に対して 教師用のセグメンテーションマップを用意することは困難であ る. セグメンテーションマップからの生成での問題を解決する ために、近年ではテキストからの画像生成が活発になっており、 簡単なテキストから画像生成が可能である上により多様な画像 を生成可能になっている.一方で,言語に依存することとユー ザが想像している画像が生成されるとは限らず何度も生成をや り直す必要がある場合があるという短所も存在する. これらの 問題を解消するために本研究では、 ラベルや言語に依存しない 方法として、図1に示すような、物体の輪郭線と色からなるカ



図 1: カラースケッチからの画像生成. (a) は物体の輪郭線を描いた線 画, (b) は (a) に対応する塗り絵を行った色画像. (c) は (a) の 下のレイヤーに (b) を置き, 合成することにより得られるカラー スケッチ画像. 画像生成モデルに (c) を入力することで (d) の フォトリアルな生成画像が得られる.

ラースケッチ画像を入力形式として用いた手法での画像生成の 精度向上を目的とする.

カラースケッチからのフォトリアルな画像の生成は、セグメ ンテーションマップからの画像生成と同様に、pix2pix [1] の手 法で実現できる. これは敵対的生成ネットワーク (Generative Adversarial Network:GAN) [2] を利用した画像変換手法の一種 である. GAN は Generator と Discriminator の 2 つのネット ワークから構成される. Generator は目的とする生成物を生成 するためのネットワークで, Discriminator は Generator が生 成した生成物と本物を識別するためのネットワークとなってい る. Generator は Discriminator を騙すように, Discriminator は Generator の生成物を見抜けるように学習を行い、2 つの ネットワークを互いに競い合わせることで本物の特徴を捉えた 生成ができるようになる.本研究では、pix2pix から派生した pix2pixHD [3] をベースとして Diffusion Model [4] で用いられ る U-Net を Generator として導入し, DPT-Large [5] を用いた 深度推定結果を活用する手法を提案し, FID 評価尺度と KID 評 価尺度,アンケートによる評価にて精度改善したことを示す.

画像生成は, GAN [2] を用いた手法が普及しており, 研究が活 発に行われている. 近年では Diffusion Model [4] を用いた手法 によりさらに高品質な生成が可能になっている. また, スケッ チやペインティングを入力とした画像生成に関する研究も提案 されている. 以下では, GAN と Diffusion Model を用いた画像 生成及びスケッチやペインティングからの画像生成に関する研 究について述べる.

# 2.1 GAN を使用した画像生成

通常の GAN はノイズのみを入力とするが, GAN から派生 した手法として条件を入力とする Conditional GAN [6] が存在 する. これによって, 条件に対応した生成ができる. Isola らの pix2pix [1] はその条件として画像を用いる手法である. pix2pix は、通常の Conditional GAN の損失関数に加えて Discriminator の損失関数に L1 損失関数が加えられ, さらに画像を小さ なパッチに分解し、パッチ単位で本物か偽物かの判別を行って いる. これにより画像の全体像を L1 損失関数で捉え, 詳細な 部分は Conditional GAN で捉えられるようになり, 精度の高 い画像変換が可能となっている.また, Generator に対しても U-Net 構造を導入する工夫を行い, さらに精度を向上させてい る. pix2pix は画像のペアがあればその対応関係を学習するこ とができるため、線画の写真化やセグメンテーションマップか らの画像生成, セマンティックセグメンテーション, モノクロ画 像のカラー化, 航空写真から地図を生成するなど様々なタスク に使用することができる.しかし、低解像度なものまでしか扱 うことができないという問題がある.

これに対して Wang らは, 高解像度な画像変換を可能にする モデルとして pix2pixHD [3] を提案した. pix2pixHD は, 高解 像度画像を効率よく判別するために異なるスケールの複数の Discriminater を用意して小さなスケールに変換した画像に対 しても Discriminater を適応できるように工夫されている. こ れにより pix2pixHD は, 2048x1024 までのより高解像度の画 像生成を可能とした. さらに, 通常の損失関数に加えて, Feature Matching Loss を定義している. これは, Discriminater に Ground Truth 画像と生成画像を入力した際で各層の出力を一 致させるためのものであり, これにより Generator により自然 な生成を促すことができる.

# 2.2 Diffusion Model を使用した画像生成

画像生成は、GAN をベースとする手法が長らく主流となって いたが、Hoらは、確率拡散モデルを用いた高品質な画像生成手法 である DDPM(Denoising Diffusion Probabilistic Models) [4] を提案した.この手法は、画像に対して小さなノイズを少しずつ 足していくことでノイズ画像に変換する過程を考え、その逆変 換によってノイズ画像から画像を生成する方法である.より具 体的には、ノイズ画像に対してかかっているノイズをニューラル ネットワークで予測し、それを繰り返すことで少しずつノイズ を取り除くことで画像を生成を行う.Diffusion Model は学習 を二乗誤差最小化で定式化でき, GAN のように Discriminator を用いたりしないことから GAN と比較すると複雑ではなく, GAN で発生しやすい勾配消失やモード崩壊などの問題の影響 を受けないという利点がある.また, GAN と比較して多様な データの生成に強いことも特徴として挙げられる [7].しかし繰 り返しノイズを予測することから GAN と比較して生成時間が かかる [8].さらに, 高解像度の画像を生成する場合には GAN と比較してより巨大なネットワークが必要となるため, モデル パラメータ数と計算量が増大する問題がある [7].

Diffusion Model では、手法そのものの工夫以外にもネット ワークアーキテクチャの工夫が行われている. 多くの Diffusion Model は DDPM [4] で有効性が示された U-Net アーキテクチャ をベースとして改良されている. U-Net [9] は, 医療用画像のセ マンティックセグメンテーション向けに提案されたスキップ接続 を特徴とする, Encoder-Decoder ネットワークである. DDPM ではこれに対して, 工夫を施した U-Net を提案している. 具体 的には, Diffusion Model として用いるために時点情報を表す 埋め込み表現の導入, Wide ResNet [10] の導入, Attention の 導入, Group Normalization [11] の導入がされている. これら の工夫により Diffusion Model での画像生成が実現された.こ れに対して Dhariwal らは、この U-Net アーキテクチャをベー スとしてアーキテクチャ探索を行い, 高品質かつ高解像度の画 像生成を実現した ADM(Ablated Diffusion Model) を提案し た[7]. 具体的にはモデルの深さ、Attention 機構のヘッド数、 ResNet Block に対して対照実験が行われ、モデルの探索が行わ れた. この研究により, Diffusion Model が当時の最先端であっ た BigGAN [12] の精度を上回ることが示された.

#### 2.3 スケッチやペインティングからの画像生成

近年、スケッチやペインティングを入力として画像生成を 実現する手法も提案されている. Isola らの pix2pix [1] では, Edges to Photo として輪郭線からのフォトリアルな画像生成を 提案した. Edges to Photo は HED (Holistically-Nested Edge Detection) [13] を用いて生成した輪郭線画像からフォトリアル な画像に復元するタスクで靴や鞄を輪郭線画像から生成する 試みが行われた. Wang らは PITI(pretraining-based imageto-image translation) [14] を用いた画像変換タスクの一つと して Sketch-to-image synthesis を提案した. Sketch-to-image synthesis についても HED [13] を用いて抽出したスケッチ画 像をフォトリアルに変換するというものである. Meng らは ペインティングからの画像生成手法として SDEdit(Stochastic Differential Editing) [15] を提案した. この手法では, 画像に 対してメディアンフィルタを適応後,減色することで画像と対 応するペインティング画像を生成し、この画像を入力として Diffusion Model によるペインティングからフォトリアルな画 像への変換を評価している. Zhang らは, 学習済みの Diffusion Model の一つである Stable Diffusion [16] に対して条件を追加 して制御するネットワーク構造である, ControlNet [17] を提案 した. この手法により, Stable Diffusion に対して様々な画像条 件付きの生成が可能となった.



図 2: 提案手法の流れ.スケッチ生成部により,写真に対応する人が描いたようなスケッチを自動生成する.Conditional GAN 部ではスケッチから 写真の状態に復元することを目的として学習を行う.

# 3 提案手法

提案手法の流れを図2に示す.ユーザは「生成したい画像の 輪郭線を描いた線画」と「線画に対応する色画像」を描き,線 画の下レイヤーに色画像を配置して合成したものを入力とする. この線画と色画像は人が描くことを想定しているが,学習時に 写真と対応する人が描いたスケッチを大量に用意することは困 難である.そこで本研究では線画と色画像を写真に対して画像 処理を行うことで自動的に抽出する.さらに,得られた線画と色 画像を乗算合成することで1枚のスケッチ画像とし,それを入 力として,Conditional GAN の学習を行う.以下では,スケッ チの生成方法及び提案する生成モデルについて説明する.

# 3.1 スケッチの生成

本研究では, スケッチの生成方法としてスケッチ1とスケッ チ2の2つの手法を用いる.

スケッチ1では,線の抽出方法としてはノイズ除去を施し,グ レースケール画像と白部分を膨張させたグレースケール画像の 差を取り白黒反転させ,輪郭線を抽出する.さらに,その画像に 対して Zhang-Suen アルゴリズム [18] を適応し細線化を行う. 色の抽出方法は,スーパーピクセルのアルゴリズムの SEEDS (Superpixels Extracted via Energy-Driven Sampling) [19] を 用いて抽出を行う.スーパーピクセルとは色が類似する領域を グルーピングするアルゴリズムで各領域を平均した色で塗りつ ぶしたものを色画像とする.

スケッチ 2 では,線の抽出方法として HED [13] を用いる. pytorch-hed [20] を用いて輪郭線を抽出し, Zhang-Suen アルゴ リズム [18] を適応し細線化を行うことで,輪郭線を抽出する. 色の抽出方法には, SDEdit [15] の Stroke Painting と同様の手 法を用いる.具体的には,カーネルサイズ 23 のメディアンフィ ルタをかけた後に Pillow [21] の quantize を用いて,6 色に減色 したものを色画像とする.

# 3.2 Conditional GAN

本研究では、カラースケッチとカラースケッチから学習済み モデルを用いて得られた深度マップを入力として pix2pixHD の Generator 部に変更を加えたモデルを提案する. 具体的には Generator を Diffusion Model に基づく U-Net に変更し, 深度 マップは DPT-Large を用いて推定する. それ以外の Discriminator 等の構造については pix2pixHD と同様である. 以下では Diffusion Model に基づく U-Net と DPT について説明する.

# 3.2.1 Diffusion Model に基づく U-Net

2.2 節で述べたように多くの Diffusion Model では DDPM [4] で有効性が示された U-Net がベースとして使用されている.本研 究では、アーキテクチャ探索により高精度化が行われた ADM [7] で使用されている U-Net を用いる. しかし, GAN では Diffusion Model のための時点情報を扱う機能を必要としないこと と、モデルパラメータ数が約 5.5 億 (554M) と非常に巨大であ ることから学習に膨大な時間を必要とし、そのまま使用するこ とは出来ない. そこで時点情報を扱う機能の削除とモデルパラ メータ数の削減を行い、pix2pixHDのGeneratorとして導入す る.具体的には、5回のダウンサンプリング数を3回まで減ら し最小の特徴量マップのサイズを 1/32\*1/32 から 1/8\*1/8 の サイズまでのダウンサンプリングまでに削減,特徴量チャンネ ル数をボトルネック部を除き削減,ボトルネック部と1/8\*1/8, 1/16\*1/16, 1/32\*1/32 の特徴量マップの領域で ResNet Block の後に配置されていた Attention Block をボトルネック部の みに削減,時点情報の埋め込み表現を計算する機構と ResNet Block から時点情報を扱う機構の削除を行った. これらの操作 によってパラメータ数を約 5000 万 (50.5M) まで削減した. 図 3 にパラメータ削減と時点情報機能の削除を行った Generator のネットワーク図を示す。

#### 3.2.2 DPT(Dense Prediction Transformers)

DPT [5] とは Ranftl らが提案した, Vision Transformer(ViT) [22] を使用したセグメンテーションモデルである.画像のセマ ンティックセグメンテーションと単眼深度推定において, CNN をベースとするモデルと比較して高精度な推定を実現させてい る.特に単眼深度推定において, DPT は物体の境界でより詳細 な出力を生成する傾向があることが示されている [5].本研究 では, ViT-Large をベースとする DPT-Large の学習済みモデ ルを用いてスケッチから深度推定を行う.そして得られた深度 マップをカラースケッチとともに GAN の入力とする.



図 3: Diffusion Model に基づく U-Net (提案手法 2 および提案手法 3 の Generator)

3.2.3 ベースライン手法

本研究ではベースライン手法として, pix2pix [1] を用いる手 法と pix2pixHD [3] を用いる手法の 2 種類を使用する. これ らのベースラインモデルを前者はベースライン手法 1,後者 はベースライン手法 2 と呼称する. ベースライン手法 2 では, pix2pixHD の Global generator network のみを Generator と して使用する.

#### 3.2.4 提案手法

本研究では提案手法として、ベースライン手法2をベース として変更を行った3つのモデルを提案する.提案手法1は、 pix2pixHD に対して DPT-Large により推定した深度マップ をカラースケッチに加えて追加の入力とするモデルであり、 pix2pixHD+DPT-Large と呼称する.提案手法2は、図3に 示した Diffusion Model に基づく U-Net を Generator とした pix2pixHD であり、pix2pixHD-Unet と呼称する.提案手法3 は提案手法1と提案手法2を組み合わせ、提案手法2に対して DPT-Large により推定した深度マップを追加の入力としたモ デルであり、pix2pixHD-Unet+DPT-Large と呼称する.

# 4 実 験

#### 4.1 データセット

本研究では, ADE20K データセット [23] と LHQ データセット [24] の 2 つを使用する.

# 4.1.1 ADE20K データセット

ADE20K データセットは屋内および屋外の様々なシーンを 含む画像データデットである. データセットは 20,210 枚の学習 用画像と 2,000 枚の検証用画像から構成される. 画像は学習用 と検証用共に 256 × 256 にリサイズし, 実験を行った.

#### 4.1.2 LHQ データセット

LHQ データセットは自然風景の画像から構成されるデータ セットである. データセットは 90,000 枚から構成される. 画像 サイズは 256 × 256 である. 本研究では, 学習用 81,000 枚, 検 証用 9,000 枚に分割し, 実験を行った.

### 4.2 評価指標

評価指標には Fréchet Inception Distance(FID) [25] と Kernel Inception Distance(KID) [26] を使用して, 生成結果の分布 と Ground truth の分布の間の特徴距離を測定する. FID スコ アと KID スコアの算出には clean-fid [27] を用いた. また, 各 モデルのパラメータ数についても評価する. パラメータ数は Generator のパラメータ数を算出し, DPT-Large 使用時はその パラメータ数も加えて評価する.

#### 4.3 アンケートによる評価

本研究では,特徴距離の評価に加えてアンケートによる人に よる評価も実施した.アンケートでは,10枚の手描きのカラー スケッチに対して各モデルを用いて生成を行い,「写実的だと 思う順番に順位づけをしてください。」と説明し,各手法の生成 結果を順位付けする方法を用いた.さらに,アンケート結果か ら各手法の票が各順位までに含まれている割合を算出した.

#### 4.4 モデルパラメータ

バッチサイズは 8, epoch 数は ADE20K データセット学習時 は 50, LHQ データセット学習時は 15 とした.また,学習率は 0.0001 としている.その他のパラメータは pix2pix はオリジナ ルと同様の値, pix2pixHD をベースとする手法は pix2pixHD と同様の値で実験を行った.

表 1: 各モデルの実験結果と各モデルのパラメータ数. FID スコアと KID スコアは生成結果と Groud Truth 間の特徴量距離を表す. パラメータ 数については Generator のパラメータ数である. ただし, DPT-Large を使用する手法についてそのパラメータ数を加えて表示している.

					スケッチ 1	$256\ge 25$	6	スケッチ 2 256 x 256			
				ADE20K		LHQ		ADE20K		LHQ	
Method			$\operatorname{Params}(M)$	FID $\downarrow$	KID $\downarrow$						
pix2pix	(ベースき	ライン手法 1)	54.40	48.18	0.0202	31.67	0.0148	116.79	0.0715	62.05	0.0396
pix2pixHD	(ベースラ	ライン手法 2)	45.59	27.74	0.0055	23.44	0.0110	53.56	0.0182	30.87	0.0140
pix2pixHD + DPT-L	arge	(提案手法 1)	389.05	33.01	0.0077	26.34	0.0132	57.93	0.0212	36.11	0.0176
pix2pixHD-Unet		(提案手法 2)	55.00	17.65	0.0013	8.56	0.0026	29.96	0.0047	9.77	0.0020
pix2pixHD-Unet + D	398.46	16.66	0.0009	6.56	0.0018	29.35	0.0046	8.89	0.0018		

			表 2:	アンケ	ートにお	ける各種	法の票数	[			
			票数				Ratio@N(N 位以上の表の割合)				
Method			1位	2位	3位	4位	5 位	Ratio@1	Ratio@2	Ratio@3	Ratio@4
pix2pix	(ベー,	スライン手法 1)	1	1	16	12	70	0.01	0.02	0.18	0.30
pix2pixHD	(ベー,	スライン手法 2)	11	16	33	32	8	0.11	0.27	0.60	0.92
pix2pixHD + DPT-L	arge	(提案手法 1)	16	14	28	37	5	0.16	0.30	0.58	0.95
pix2pixHD-Unet		(提案手法 2)	36	34	6	11	13	0.36	0.70	0.76	0.87
pix2pixHD-Unet + DPT-Large (提案手法 3)			36	35	12	15	2	0.36	0.71	0.83	0.98

# 4.5 実験結果

表1はADE20K データセットとLHQ データセットでの実験 結果と各モデルのパラメータ数である.表1より,pix2pixHD のGeneratorをDiffusion Modelに基づくU-Netした提案手法 2はpix2pixHDと比較して両方のデータセットの各スケッチ生 成方法で精度が向上していることがわかる.さらにDPT-Large を用い深度マップを追加の入力とした提案手法3はそれぞれの データセットで最も高い精度となった.しかし,pix2pixHDに 対して深度マップを追加の入力として与えた提案手法1に関し てはpix2pixHDに対して精度が低下した.

図4はLHQ データセットのスケッチ1での生成例である. 図4より pix2pix では細かな質感の再現が上手くできておらず、 ぼやけた印象となっていることが確認できる. pix2pixHD は, 質感については pix2pix よりも良いものの輪郭がはっきりしな い画像となっていることが確認できる.提案手法についてはこ れらの問題が軽減されており、質感や輪郭のぼやけが改善され ていることが確認できる.しかし,建物が含まれる例については 提案手法2と提案手法3においても余計な質感が付与されてい たりする例が確認できた. 図5の LHQ データセットのスケッ チ2での生成例については、pix2pixでは図4での生成例と同 様に質感がのっぺりとしていることが確認でき, pix2pixHD に ついても輪郭のぼやけが確認できる. これに対して提案手法で は物体の輪郭がはっきりとし、細かな質感が生成されているこ とが確認できる. さらに, Ground truth と比較しても遜色ない 品質で復元できていることが確認できる.しかし、人工物が含 まれる例については上手く生成できない箇所が含まれる例も確 認できた.

また,図6はペイントソフトを用いて人が描いたスケッチからの生成結果であり,LHQデータセットでスケッチ2での学習を行った各モデルを用いて生成を行った.図6より,人が描いたスケッチからの生成においてもフォトリアルな画像を生成で

きていることが確認できる.特に,提案手法2と提案手法3は 草木や岩肌の質感や水面の表現が上手く生成できていることが 確認できる.しかし,図4の一番下の例のように,空部分に不自 然なテクスチャが現れるケースも確認できた.

4.5.1 アンケート結果

図 6 の結果を用いて 10 名に対してアンケートを実施した結 果が, 表 2 である. Ratio@N は全体の票に対して N 位以上に選 ばれた割合を表す. 表 2 より提案手法 3 は最も上位に選ばれる 票が多く, 下位に票が入ることが少ないことが分かる. ただし, 1 位に選ばれた数は提案手法 2 と提案手法 3 で同数の 36 票で あることが確認できる.

#### 4.6 考 察

実験結果より, Diffusion Model に基づく U-Net を導入した 提案手法2は pix2pix 程度のパラメータ数を保ちながら高精度 化することに成功した. これにより, Diffusion Model に基づく U-Net が Diffusion Model だけでなく GAN の Generator と しても効果的であることが確認できた.これは、Wide ResNet や Attention, Group Normalization が GAN に対しても有効 であったことが考えられる. また, pix2pixHD の Generator が ダウンサンプリング後のボトルネック部に ResNet が配置され ているのに対し, Diffusion Model に基づく U-Net ではダウン サンプリング毎に ResNet が配置されており, 大きい解像度の特 徴量マップに対しても ResNet を適応させていることが精度の 向上につながっていると考えられる. さらに, DPT-Large を用 いて深度マップを追加で与える拡張を行った提案手法3は,提 案手法2よりさらに精度を改善した. これは、深度マップにより 同じ物体の領域が近い深度でグルーピングされたことが, 質感 の描き分けや空間表現の向上につながったと考えられる.しか し, pix2pixHD に対して深度マップを追加の入力として与えた 提案手法1については、精度が低下する結果となった.これは、



図 4: LHQ データセットでのスケッチ 1 の各手法の生成例. Ground truth は実際の写真, Color Sketch は入力するカラースケッチである.



図 5: LHQ データセットでのスケッチ 2 の各手法の生成例. Ground truth は実際の写真, Color Sketch は入力するカラースケッチである.



図 6: LHQ データセットでのスケッチ 2 の人が描いたカラースケッチからの生成例. Color Sketch は入力となるペイントソフトを用いて人が描いたカラースケッチある.

深度マップが質感の向上に影響するほどベースの pix2pixHD の精度が良くなかったためであると考えられる.

また, 各手法の LHQ データセットでの FID スコアと KID ス コアは, ともに ADE20K データセットでの実験結果よりも高 精度な結果が得られた. これは, ADE20K データセットが様々 なシーンを含むデータセットであることに対して, LHQ データ セットが自然風景のみのデータセットであり, ADE20K データ セットのほう変換が困難であったためだと考えられる. さらに, LHQ データセットでの結果で人工物が上手く生成できない例 が得られたが, これは, LHQ データセットに人工物を含む画像 が少なかったためだと考えられる.

人が描いたスケッチからの生成結果については、人が描いた スケッチからであってもフォトリアルに変換可能であることが 確認できたが、自動生成したスケッチからの生成と比較すると 細かな質感の面で劣っていることが確認できた.これは、自動 生成により得られたスケッチのみを学習の入力としていたため に, スケッチのスタイルに制約が生まれてしまっているためだ と考えられる.

# 5 おわりに

本研究では、ラベルや言語に依存しない画像生成の入力形 式として物体の輪郭線と色からなるカラースケッチ画像を入 力とする画像生成において、生成手法として Diffusion Model に基づく U-Net と DPT-Large を導入した pix2pixHD を提 案した. 実験結果より、Diffusion Model に基づく U-Net と DPT-Large による深度推定結果の活用はスケッチからの画 像生成において精度の向上に効果的であることが確認でき た. 本研究では、DPT-Large による深度推定結果の活用が精 度の向上につながったが、画像のセグメンテーションは深度 推定だけでなく、セマンティックセグメンテーションが存在 し、STEGO(Self-supervised Transformer with Energy-based Graph Optimization) [28] や PiCIE(Pixel-level feature Clustering using Invariance and Equivariance) [29] のような教師 なしセマンティックセグメンテーションモデルの活用について も検討する必要があると考えられる.

また, 人が描いたスケッチから生成が可能であることも確認 したが, 自動生成したカラースケッチからの生成結果と比較す ると細かな質感の面で劣っており, より多様なスタイルのスケッ チや粗いスケッチには対応することが出来ないと考えられる. また, ラベルに依存しない入力形式としてスケッチ画像を提案 したが, 現在の入力画像は Ground truth 画像の情報を多く保 持しており, 入力画像の時点である程度フォトリアルであると 言える. Ground truth 画像と比較すると粗いが, より粗くスタ イルが異なっていてもフォトリアルに変換可能な手法が必要で あると考えられる.

今後の課題としては, Diffusion Model を使用し SDEdit [15] を導入することでスケッチのスタイルへの依存を軽減すること や ControlNet を導入し, 学習済み Diffusion Model を用いる ことによる精度改善が考えられる.

# 文 献

- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017.
- [2] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- [3] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018.
- [4] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, Advances in Neural Information Processing Systems, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.
- [5] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (*ICCV*), pages 12179–12188, October 2021.
- [6] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets, 2014.
- [7] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 8780–8794. Curran Associates, Inc., 2021.
- [8] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference* on *Learning Representations*, 2021.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. Unet: Convolutional networks for biomedical image segmentation. CoRR, abs/1505.04597, 2015.
- [10] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *BMVC*, 2016.
- [11] Yuxin Wu and Kaiming He. Group normalization. In Proceedings of the European Conference on Computer Vision (ECCV), September 2018.
- [12] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis.

In International Conference on Learning Representations, 2019.

- [13] Saining Xie and Zhuowen Tu. Holistically-nested edge detection. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), December 2015.
- [14] Tengfei Wang, Ting Zhang, Bo Zhang, Hao Ouyang, Dong Chen, Qifeng Chen, and Fang Wen. Pretraining is all you need for image-to-image translation. In arXiv, 2022.
- [15] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2022.
- [16] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.
- [17] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023.
- [18] C. Y. Suen T. Y. Zhang. A fast parallel algorithm for thinning digital patterns. In *Communications of the ACM*, volume 27, pages 236–239, 1984.
- [19] Michael Van den Bergh, Xavier Boix, Gemma Roig, and Luc Van Gool. SEEDS: superpixels extracted via energydriven sampling. *CoRR*, abs/1309.3848, 2013.
- [20] Simon Niklaus. A reimplementation of HED using PyTorch, 2018.
- [21] Clark Alex. Pillow (pil fork) documentation, 2015.
- [22] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In International Conference on Learning Representations, 2021.
- [23] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017.
- [24] Ivan Skorokhodov, Grigorii Sotnikov, and Mohamed Elhoseiny. Aligning latent and image spaces to connect the unconnectable. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14144–14153, October 2021.
- [25] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017.
- [26] Mikołaj Bińkowski, Dougal J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD GANs. In International Conference on Learning Representations, 2018.
- [27] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties in gan evaluation. In CVPR, 2022.
- [28] Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snavely, and William T. Freeman. Unsupervised semantic segmentation by distilling feature correspondences. In International Conference on Learning Representations, 2022.
- [29] Jang Hyun Cho, Utkarsh Mall, Kavita Bala, and Bharath Hariharan. Picie: Unsupervised semantic segmentation using invariance and equivariance in clustering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 16794–16804, June 2021.