

線画のフレーム補間を用いたアニメ自動着色

陸 儀^{1,a)} 中島 克人^{1,b)}

概要：着色済みのアニメ画像の着色情報を伝播させる事による後続フレームのアニメ線画群の自動着色を目指している。これまでの研究で、着色対象のアニメ線画に擬似陰影を付与する事で精度の高い着色が出来る事を示したが、フレーム間でアニメキャラクターの移動などによるコンテンツの変動が大きい場合に色伝播の精度が落ちるという課題があった。そこで、線画フレーム間に深層学習を用いたフレーム補間技術によって補助的な中間フレームを設ける事により、フレーム間の内容の変動量を減少させ、色伝播の精度向上と着色品質の改善を行ったので、その評価結果を報告する。

1. はじめに

日本のアニメ産業では、市場拡大に伴う制作タイトル数の増加と作画品質の要求の高度化によりアニメ制作会社の負担が増大しており、情報技術を用いたアニメ制作の労力削減が検討されている。本稿ではアニメ制作において多くの人手を要する着色工程の支援手法を提案する。

アニメの着色工程では、事前に定められた色彩設定に従ってモノクロの動画に色を付ける必要がある。そのため完全自動着色ではなく、あらかじめ動画のキーフレームを手で着色した後に、キーフレームと後続フレームの間でフレーム間追跡を行い、キーフレームの色を後続フレームに伝播する着色手法 [1-5] が用いられる。しかしアニメの線画は輝度勾配を含まず、フレーム間でアニメキャラクターの移動などによるコンテンツの変動が大きいと、フレーム間追跡の精度が低くなるという課題がある。陸ら [4] は入力線画に対し陰影推定モデルを用いて輝度勾配の代わりとなる擬似的な陰影を付与する手法を用い、フレーム間でコンテンツ内の各部の対応付けを容易にする手法を提案した。この手法はアニメ線画の自動着色の品質をある程度向上させたが、フレーム間でコンテンツの変動量大きい場合にはコンテンツ各部の対応付けが困難となり、追跡精度が低下するという課題が示されていた。

本稿ではこの課題に対処するために、本来のフレーム間に補助的な補間フレームを生成して線画フレーム間のコンテンツの変動量を減少させる手法を提案する。

提案手法の着色例を図 1 に示す。フレーム補間手法に

は DAIN [6] と CAIN [7] の 2 種類を用い、それぞれでフレーム補間したデータを用いてフレーム間追跡モデルの MAST [8] に学習させ、最終的な着色精度を比較する。

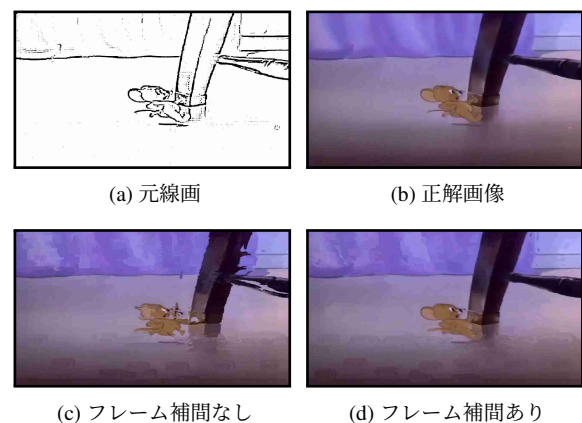


図 1 提案手法の着色結果の例

2. 関連研究

2.1 自動着色のためのフレーム間追跡

フレーム間追跡には、物体単位の追跡と画素単位の追跡の 2 つの方法がある。物体単位の追跡 (object tracking) は、2 枚の入力画像間で最も類似度の高い物体同士を組として出力する。画素単位の小さな対応の失敗が生じにくいという長所があるが、領域分割の精度が着色品質に影響を与えることと領域内部の画素同士を対応付けられずグラデーションを上手く伝播できないことが短所である。深層モデルを用いる手法には RVOS [9]、非深層モデルによる手法には Ramos ら [1]、Maejima ら [3]、陸ら [5] の手法がある。

画素単位の追跡 (dense tracking) は、2 枚の入力画像の画素間の類似度を求めて類似度行列として出力する。類似度

¹ 東京電機大学大学院
Graduate School of Tokyo Denki University
a) 20fmi24@ms.dendai.ac.jp
b) nakajima@mail.dendai.ac.jp

行列を用いることでグラデーションも伝播できる長所があるが、1画素～数画素単位の細かな対応失敗が生じやすいという短所がある。深層モデルによる手法には、MAST [8], MuG [10], CorrFlow [11], CycleTime [12], Tracking via Colorization [13] がある。しかし、提案されているいずれの手法も、フレーム間でコンテンツの変動が大きい場合には追跡が困難になる。

2.2 動画の情報補間によるフレーム間追跡の精度向上

陸ら [4] は、アニメーション線画に対して陰影推定モデルを用いて擬似陰影を付与することで、自動着色の精度が向上することを示した。陸らは陰影推定モデルに ShadeSketch [14] と Pix2PixHD [15] を用いた場合の着色精度を比較した。ShadeSketch による擬似陰影の例を図 2 に示す。



図 2 ShadeSketch による擬似陰影の例

Narita ら [16] はアニメーション線画に対して距離変換を用いて勾配を付与することで、線画のオプティカルフロー計算の精度が向上することを示した。しかし、これらの手法もフレーム間のコンテンツ変動が大きい場合にはフレーム間のコンテンツ対応は困難である。

2.3 フレーム補間

フレーム補間とは、時系列的に前後関係のある 2 枚の入力フレーム間でフレーム内コンテンツの移動軌跡を推定

し、時系列の中間に位置するフレームを生成することである。生成された中間フレームを用いて再度フレーム補間を繰り返すことで、2, 4, 8, ..., 2^n 倍に補間することが可能である。深層学習に基づくフレーム補間モデルには DAIN [6] と CAIN [7] がある。DAIN はフレーム補間の計算時にオプティカルフローに加えてコンテンツの深度情報を推定することで、物体の一部が隠れている場合にも高精度なフレーム補間を実現する。CAIN は従来のフレーム補間で用いられてきたオプティカルフロー計算をアテンション機構に基づく特徴変換操作に置き換えることで、高速かつ省メモリなフレーム補間を実現する。CAIN は DAIN に比べて実行時間を約 14 分の 1 に削減した。

3. 提案手法

入力線画に対しフレーム補間を適用した後に、フレーム間追跡により参照フレームの色を後続フレームに伝播するアニメ自動着色手法を提案する。提案手法の処理フローを図 3 に示す。提案手法の入力データは参照フレーム(カラー)、参照フレーム(線画)、着色対象フレーム(線画)の 3 種類である。参照フレームとは着色対象フレームの一部を人手で着色したフレームであり、参照フレーム群には少なくともアニメを構成する各カットの先頭フレームを含む必要がある。

提案手法は大きく分けてフレーム補間、画素単位のフレーム間追跡、色伝播の 3 つの工程で構成される。フレーム補間とは、時系列的に前後関係のある 2 枚の入力フレームから中間フレームを出力する工程である。画素単位のフレーム間追跡とは、2 枚の入力フレームに含まれる画素間の類似度を実数値として求め、類似度を縦横の行列形式で格納した類似度行列を出力する工程である。色伝播とは、2 枚の入力フレームの類似度行列に基づき、一方のカラーフレームの画素値をもう一方の白黒またはグレースケールのフレームで最も類似度の高い画素にコピーする工程である。

3.1 フレーム補間

本稿ではフレーム補間モデルとして DAIN [6] と CAIN [7] の 2 種類のモデルを比較する。DAIN と CAIN のフレーム補間には、計算効率と補間品質の 2 点の違いがある。計算効率に関しては $1,280 \times 720$ の画像 1 枚あたり、DAIN は約 3,500ms、CAIN は約 250ms でフレーム補間できる。DAIN と CAIN のフレーム補間の例を図 4 に示す。図 4 より補間品質に関しては、DAIN の方が CAIN よりも鮮明な輪郭線を持つ中間フレームを生成する傾向がある。輪郭線の傾向の違いが自動着色に与える影響は本稿の評価実験で検証する。

3.2 画素単位のフレーム間追跡

画素単位のフレーム間追跡には、MAST [8] を用いる。

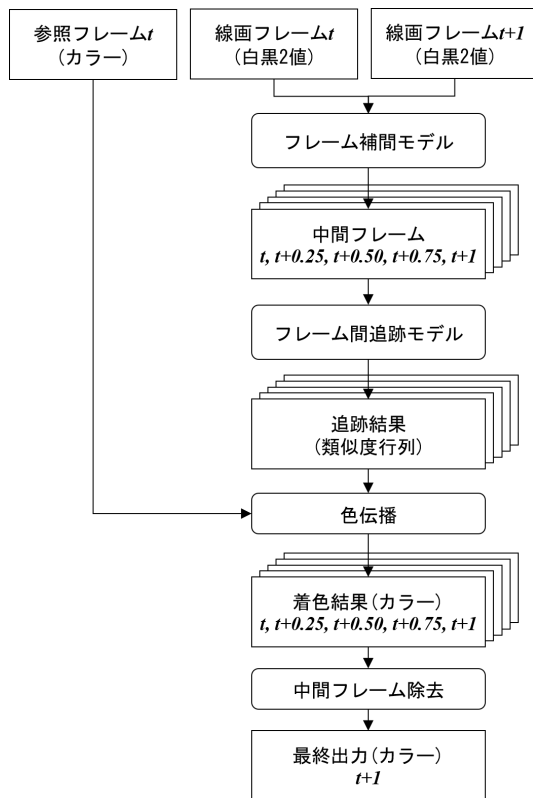


図3 提案手法の処理フロー

MASTの特徴は二つある。一つは自己教師学習であることで、もう一つは追跡範囲の局所化により計算効率の向上を図ったことである。ここで、自己教師学習とは、訓練データとしてフレーム間の画素の対応関係のアノテーションが不要で、カラー動画と対応するグレースケール動画のみでモデル学習できることを意味する。

計算効率の向上は、画素の対応関係を求める範囲を局所化して類似度行列のサイズを削減したことで、これにより高解像度のフレーム間追跡も可能にした。従来の画素単位の追跡手法は、フレーム間の全画素同士の類似度を求め類似度行列として出力していた。そのため入力画像の解像度が $W \times H$ のとき類似度行列は $(WH) \times (WH)$ の巨大な行列となり、高解像度の動画を処理するには膨大な処理時間とGPUのメモリ量を必要としていた。自動着色タスクにおいては出力動画の解像度は直接に人の視野に反映される重要な項目であるため、本稿のフレーム間追跡モデルとして高解像度なフレーム間追跡が可能なMASTを用いた。

3.3 色伝播

色伝播では類似度行列に含まれる画素間の類似度の値に基づき、参照フレームの画素値を着色対象フレームで最も類似度の高い画素にコピーする。本稿では最も類似度の高い画素に単純にコピーする方法を用いた。しかし、位置関係が一定の距離の範囲内の画素という条件内で最も高い画素にコピーしたり、類似度行列に対して事前に平滑化フィ

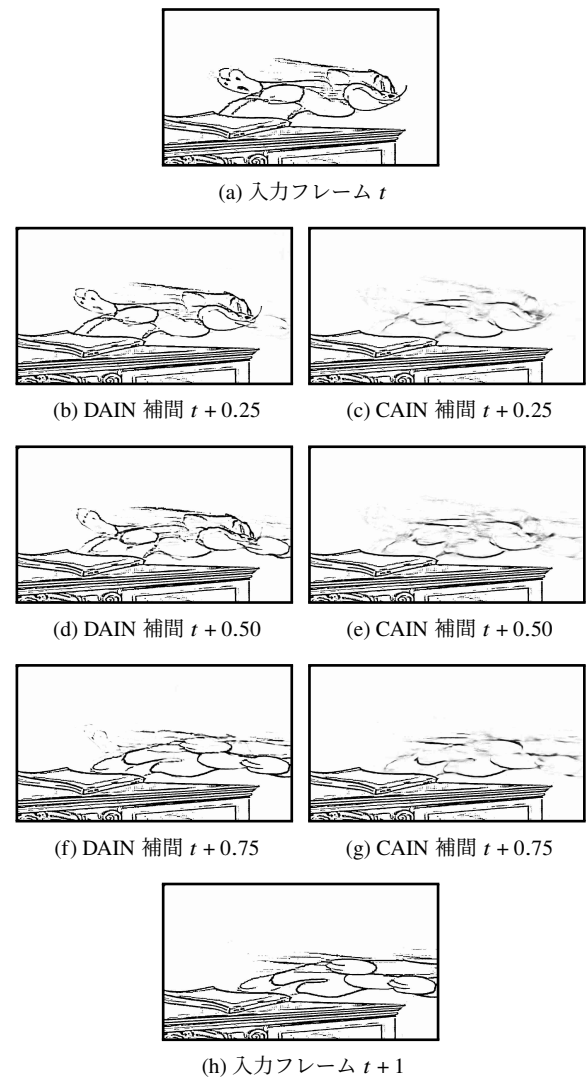


図4 DAINとCAINによるフレーム補間の例

ルタを施してフレーム間追跡で生じたノイズの影響を減らす後処理も考えられ、これらは今後の課題である。

4. 評価

フレーム補間を用いたデータセットと用いないデータセットの2種類でMASTを学習し、フレーム補間の有無がフレーム間追跡に与える影響を評価する。また、擬似陰影を用いる陸らの先行手法[4]と提案手法の精度を比較する。

4.1 学習設定

フレーム補間モデルのDAINとCAIN、および擬似陰影推定モデルのShadeSketchには著者らが公開している学習済みモデルを用いた。MASTはモデルの入力解像度 $2,560 \times 1,440$ 画素、バッチサイズ8、初期学習率0.001、最適手法Adam[17]を用いて35エポック学習した。MASTはモデルの入力解像度が $2,560 \times 1,440$ 画素のとき、出力される類似度行列のサイズは 640×360 となるため、色伝播後の着色結果の解像度は 640×360 画素である。誤差関

数にはスムーズ L1 損失 [18] を用いた。

4.2 データセット

MAST を学習するために著作権期間が満了したトムとジェリーの第 1 シリーズ [19] の動画を使用した。データセットとしては $t, t+1$ の 2 枚の連続フレームを 1 組として、訓練セットは第 5 話から 450 組、テストセットは第 6 話から 450 組をそれぞれランダムに選択した。

4.3 評価方法

フレーム補間の効果を評価するため、以下の 5 つの手法による自動着色の結果を評価した。

- (1) 実陰影 (カラー画像をグレースケール化して生成)
- (2) 元線画 (加工なし)
- (3) 擬似陰影
- (4) フレーム補間 (DAIN)
- (5) フレーム補間 (CAIN)

図 5 は (1), (2) と (3) のデータセットを用いる際の MAST の入出力を表す。(1) は比較のために用意したデータセットであり、着色対象線画に対応する人手着色の正解のカラーフレームをグレースケール化したものである。(2) は着色対象線画をそのまま入力する場合である。(3) は (2) に対して [4] による擬似陰影を付す場合である。

(4) と (5) は図 3 において、MAST によるフレーム間追跡の前にフレーム補間を行い、中間フレームも含めて MAST に入力する場合である。図 4 は DAIN および CAIN によるフレーム補間について更に詳しく示したものとなる。

なお、今回の評価に必要な着色対象線画等のためのカラーフレームの線画化には XDoG [20] を用いている。XDoG のパラメータは人手で調整し、 $(kernel, \sigma, \epsilon, \phi, K, \gamma) = (17, 1.0, -9.9, 11.1, 5.2, 0.999)$ を用いた。

4.4 評価指標

評価は図 3 や図 5 における色伝播結果の着色結果画像とこれに対応する人手着色の正解のカラー画像との比較で行い、評価指標には NRMSE と SSIM を用いた。NRMSE は画素値の差を正規化した指標である。NRMSE の導出式を式 1 に示す。

$$NRMSE(x_1, x_2) = \frac{\sqrt{\frac{\sum_{i=1}^n (x_{1,i} - x_{2,i})^2}{n}}}{x_{max} - x_{min}} \quad (1)$$

2 枚の同じ解像度の入力画像を x_1, x_2 とする。 $x_{1,i}, x_{2,i}$ はそれぞれの画像の i 番目の画素を表す。 n は画像 1 枚あたりの総画素数である。 x_{max} と x_{min} はそれぞれ $x_{1,i}, x_{2,i}$ の最大値と最小値である。NRMSE は相違度指標であるため、NRMSE の値を 1.0 から引くことで類似度指標として用いる。SSIM は 2 枚の画像の類似度を輝度、コントラスト、物体構造の観点から求めた指標である。SSIM の導出

式を式 2 に示す。

$$SSIM(x_1, x_2) = \frac{(2\mu_{x_1}\mu_{x_2} + c_1)(2\sigma_{x_1x_2} + c_2)}{(\mu_{x_1}^2 + \mu_{x_2}^2 + c_1)(\sigma_{x_1}^2 + \sigma_{x_2}^2 + c_2)} \quad (2)$$

μ_x は x の平均値、 σ_{ab} は a, b の共分散、 σ_a^2 は x の分散、 $c_1 = (k_1 \times (2^d - 1))^2$ 、 $c_2 = (k_2 \times (2^d - 1))^2$ 、 d はビット深度であり、 $k_1 = 0.01$ 、 $k_2 = 0.03$ とした。

4.5 評価結果と考察

学習終了時の精度を表 1 に、学習曲線を図 6 に示す。また着色結果例を図 7 に示す。

表 1 学習終了時の着色精度

	1-NRMSE	SSIM
(1) 実陰影	0.906	0.876
(2) 元線画	0.852	0.804
(3) 擬似陰影	0.853	0.796
(4) DAIN 補間	0.890	0.835
(5) CAIN 補間	0.889	0.830

表 1 より (4)DAIN と (5)CAIN によるフレーム補間で、(2)元線画よりも高い着色精度を得た。フレーム間でコンテンツの動きが大きい場合に着色精度が低下する問題が陸ら [4] により指摘されていたが、フレーム補間を用いることでこの問題を改善できることが確認できる。

本評価では (4) と (5) のフレーム補間が (2) 元線画に対して着色精度を改善する一方で、(3) 擬似陰影は着色精度を改善できなかった。擬似陰影の先行手法と今回は以下の点で異なり、下記の違いが擬似陰影の効果に影響を及ぼしたと考えられる。

- (1) フレーム間追跡モデルに Tracking via Colorization ではなく MAST を使用している
- (2) モデルの入出力解像度が異なる。先行手法は入力 $1,024 \times 1,024$ 画素、出力 128×128 画素、本評価では入力 $2,560 \times 1,440$ 画素、出力 640×360 画素である
なお、フレーム補間して生成したデータセットに対し、さらに擬似陰影を付与する場合としない場合を比較するため、それぞれについて学習を行い、学習曲線を比較した (図 8)。DAIN に関しては、学習初期の曲線の値の推移には差が見られるが、収束後のロスや SSIM 精度には差が見られない。CAIN に関しては、擬似陰影を付与した場合には学習中の SSIM 精度の振動が抑制できているが、学習全体を通してのモデルの最高精度は 0.01 ポイント低下した。即ち、今回の評価では、フレーム補間に対して追加で擬似陰影を組み合わせても着色精度は改善しなかった。

図 7 では DAIN は明確に優れた着色品質を実現している。図 7e の DAIN の着色結果では、椅子の脚の部分やジェリー (茶色いねずみのキャラクター) において図 7b の実

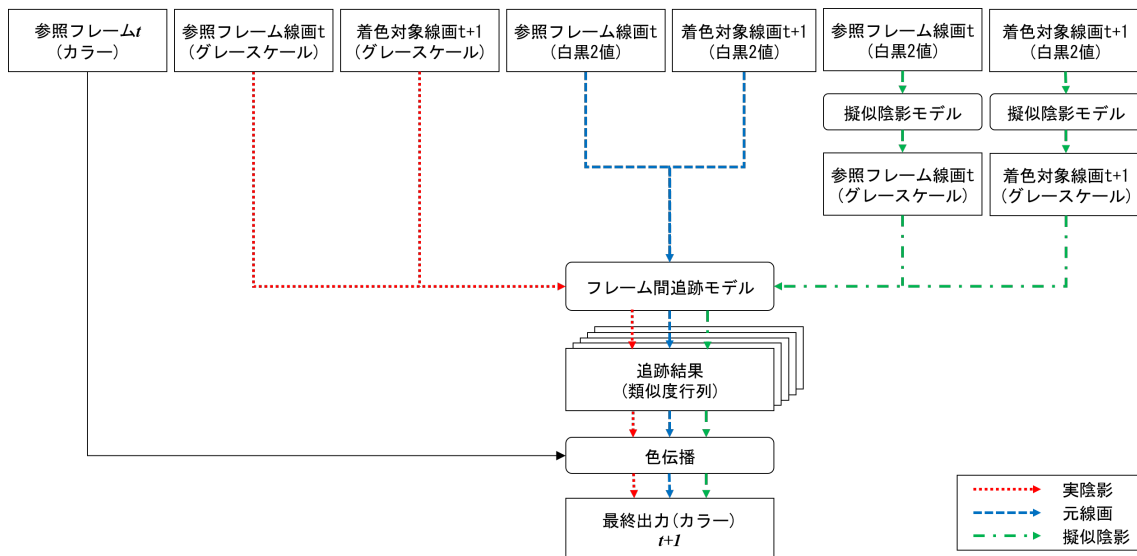


図5 比較対象データセットの処理フロー

陰影を上回る着色品質を実現している。

図7dのCAINの着色結果では、椅子の脚の輪郭部分の塗りが不規則に欠ける問題が見られた。その理由を図9に示す。図9c, dはDAINとCAINのフレーム補間結果の間に顕著な差があるフレームを示したものである。入力画像a, bはそれぞれ $t, t+1$ の線画であり、出力画像c, dはそれぞれa, bをDAINとCAINで補間した $t+0.75$ のフレームである。

なお図7, 9は同一の入力フレームを用いている。図9dのCAINの補間結果を見ると椅子の脚やジェリーの体の補間に失敗し輪郭線の一部が消えている。図7hを見ると、図9dで輪郭線が消えている部分の着色に失敗していることが確認できる。図9cのDAINには補間の失敗は見られず、図7gにおいても高品質な着色結果が確認できる。

図7ではDAINの着色品質が実陰影を上回るように見えるが、表1の客観評価ではDAINの精度は実陰影に及んでいない。またDAINとCAINの精度も大差はない。その理由はDAINはフレーム補間に失敗した場合に大きなノイズを生成することがあるためである。図10に例を示す。

図10cのDAINによるフレーム補間結果では、背景領域に対して広範囲にノイズ画素が出現しており、図10eの着色結果にもノイズが出現している。図10dのCAINによる補間結果では輪郭線が膨張する傾向は見られるが、背景にノイズ画素は出現していない。DAINは広範囲にわたって単色や緩やかなグラデーションが存在する部分にノイズを生成する傾向がある。

DAINとCAINに関する評価結果のまとめを表2に示す。DAINとCAINはどちらも十分な着色結果の改善を示し、特にDAINは実陰影を上回る品質での着色も可能なことを示した。DAINとCAINにはどちらも補間失敗時に着色品質を低下させる課題があるため、補間失敗の影響と頻度を

低くすることが着色品質の改善に有効であると考えられる。

表2 DAINとCAINの比較

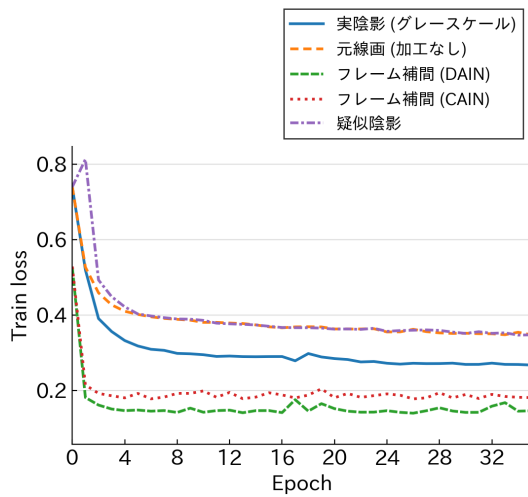
比較項目	DAIN	CAIN
着色精度の改善	大きい	大きい
補間失敗時	ノイズ出現	輪郭線消失
補間失敗の影響	大きい	小さい
補間失敗の頻度	低い	高い

5. おわりに

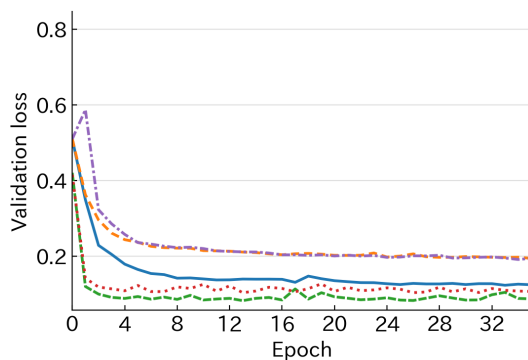
我々はアニメ線画の自動着色において、着色前に線画をフレーム補間することで着色精度が向上することを示した。フレーム補間モデルにはDAINとCAINを比較した。DAINを用いると実陰影の着色品質を上回る場合もあることが確認できた。DAINとCAINはフレーム補間に失敗した場合に着色品質を低下させる課題がある。補間失敗時の影響と頻度の減少が着色品質改善に有効であると考えられる。

参考文献

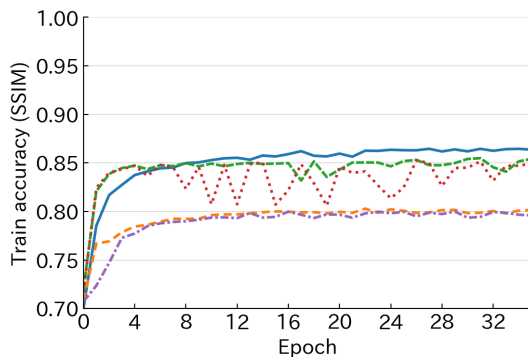
- [1] Ramos, A. and Flores, F.: Colorization of Grayscale Image Sequences using Texture Descriptors, *Proc. 14th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISIGRAPP)* (2019).
- [2] Ishii, D., Kubo, H., Shinagawa, S., Maejima, A., Funatomi, T., Nakamura, S. and Mukaigawa, Y.: Confidence-aware Practical Anime-style Colorization, *Proc. ACM SIGGRAPH 2020* (2020).
- [3] Maejima, A., Kubo, H., Funatomi, T., Yotsukura, T., Nakamura, S. and Mukaigawa, Y.: Graph Matching based Anime Colorization with Multiple References, *Proc. ACM SIGGRAPH 2019* (2019).
- [4] 陸 儀, 中島克人: 疑似陰影を用いたフレーム間追跡によるアニメ線画の自動着色, 第23回画像の認識・理解シ



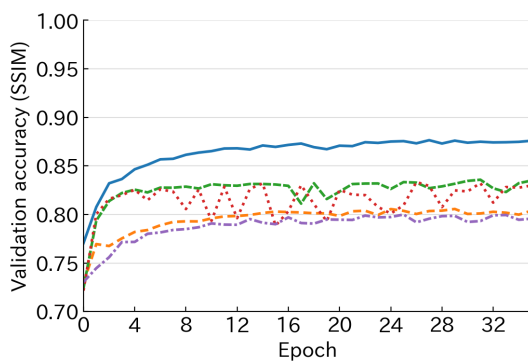
(a) ロスの推移 (訓練セット)



(b) ロスの推移 (テストセット)



(c) SSIM 精度の推移 (訓練セット)



(d) SSIM 精度の推移 (テストセット)

図 6 提案手法の学習曲線

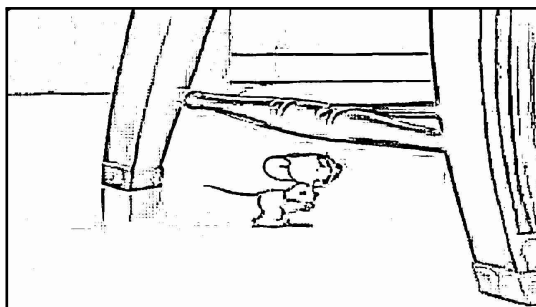
- ンボジウム (MIRU) (2020).
- [5] 陸 儀, 中島克人: 着色済み参照フレームを用いるアニメ線画の自動着色手法の提案, 情報処理学会第 82 回全国大会講演論文集 (2020).
 - [6] Bao, W., Lai, W.-S., Ma, C., Zhang, X., Gao, Z. and Yang, M.-H.: Depth-Aware Video Frame Interpolation, *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019).
 - [7] Choi, M., Kim, H., Han, B., Xu, N. and Lee, K.: Channel Attention Is All You Need for Video Frame Interpolation, *Proc. Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)* (2020).
 - [8] Lai, Z., Lu, E. and Xie, W.: MAST: A Memory-Augmented Self-Supervised Tracker, *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020).
 - [9] Ventura, C., Bellver, M., Girbau, A., Salvador, A., Marques, F. and Giro-i Nieto, X.: RVOS: End-To-End Recurrent Network for Video Object Segmentation, *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019).
 - [10] Lu, X., Wang, W., Shen, J., Tai, Y.-W., Crandall, D. and Hoi, S.: Learning Video Object Segmentation from Unlabeled Videos, *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020).
 - [11] Lai, Z., Lu, E. and Xie, W.: Self-supervised Learning for Video Correspondence Flow, *Proc. 30th Annual British Machine Vision Conference (BMVC)* (2019).
 - [12] Wang, X., J. A. and Efros, A.: Learning Correspondence From the Cycle-Consistency of Time, *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019).
 - [13] Vondrick, C., Shrivastava, A., Fathi, A., Guadarrama, A. and Murphy, K.: Tracking Emerges by Colorizing Videos, *Proc. 15th European Conference on Computer Vision (ECCV)* (2018).
 - [14] Zheng, Q., Li, Z. and Bargteil, A.: Learning to Shadow Hand-drawn Sketches, *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020).
 - [15] Wang, T.-C., Liu, M.-Y., Zhu, J.-Y., Tao, A., Kautz, J. and Catanzaro, B.: High-Resolution Image Synthesis and Semantic Manipulation with Conditional GANs, *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2018).
 - [16] Narita, R., Hirakawa, K. and Aizawa, K.: Optical Flow Based Line Drawing Frame Interpolation Using Distance Transform to Support Inbetweens, *Proc. 26th IEEE International Conference on Image Processing (ICIP)* (2019).
 - [17] Kingma, D. and Ba, J.: Adam: A Method for Stochastic Optimization, *Proc. 3rd International Conference on Learning Representations (ICLR)* (2015).
 - [18] Girshick, R.: Fast R-CNN, *Proc. IEEE International Conference on Computer Vision (ICCV)* (2015).
 - [19] Hanna, W. and Barbera, J.: Tom and Jerry Theatrical Cartoon Series, *Metro-Goldwyn-Mayer* (1940 to 1958).
 - [20] Winnemöller, H., Kyprianidis, J. and Olsen, S.: XDoG: An eXtend difference-of-Gaussians compendium including advanced image stylization, *Computers & Graphics*, Vol. 36, No. 6 (2012).



(a) 参照フレーム t



(b) 正解画像 $t+1$



(c) 着色対象線画 $t+1$



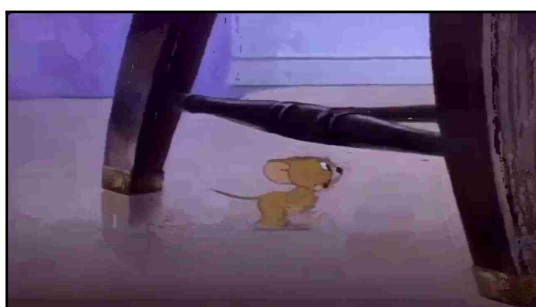
(d) 実陰影の着色結果 $t+1$



(e) 元線画の着色結果 $t+1$



(f) 擬似陰影の着色結果 $t+1$

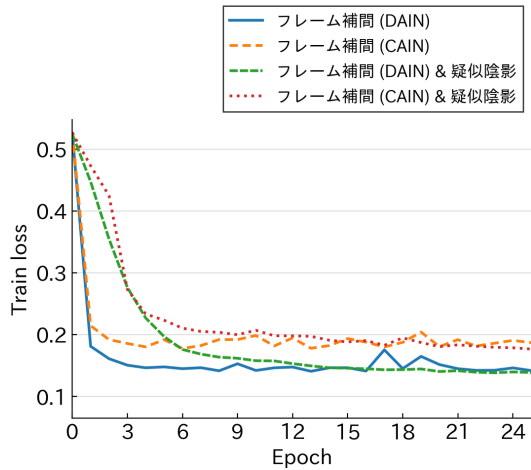


(g) DAIN の着色結果 $t+1$

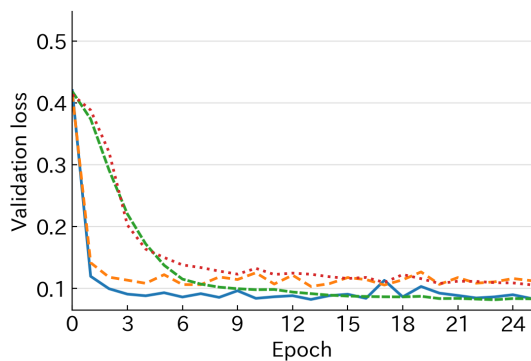


(h) CAIN の着色結果 $t+1$

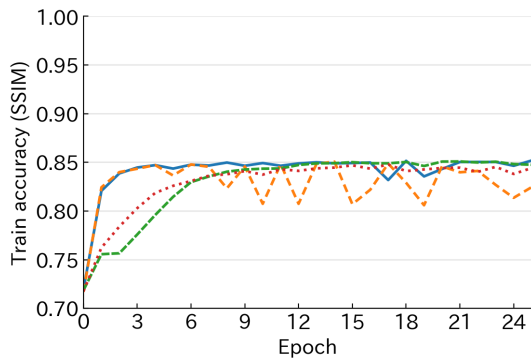
図7 着色結果の比較



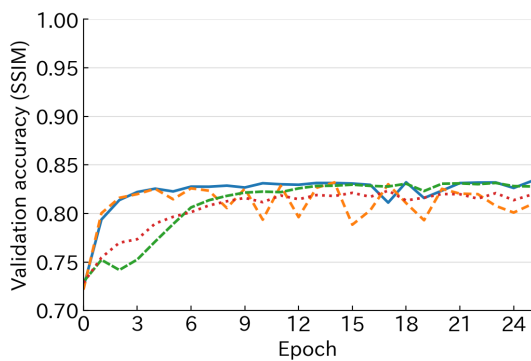
(a) ロスの推移 (訓練セット)



(b) ロスの推移 (テストセット)



(c) SSIM 精度の推移 (訓練セット)



(d) SSIM 精度の推移 (テストセット)

図8 フレーム補間と疑似陰影を組み合わせた際の学習曲線

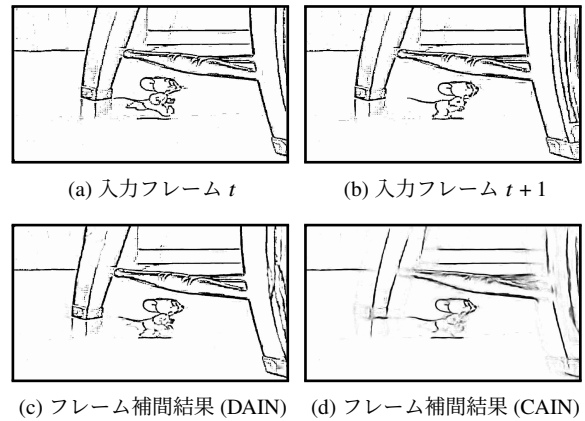


図9 フレーム補間による輪郭線消失の例

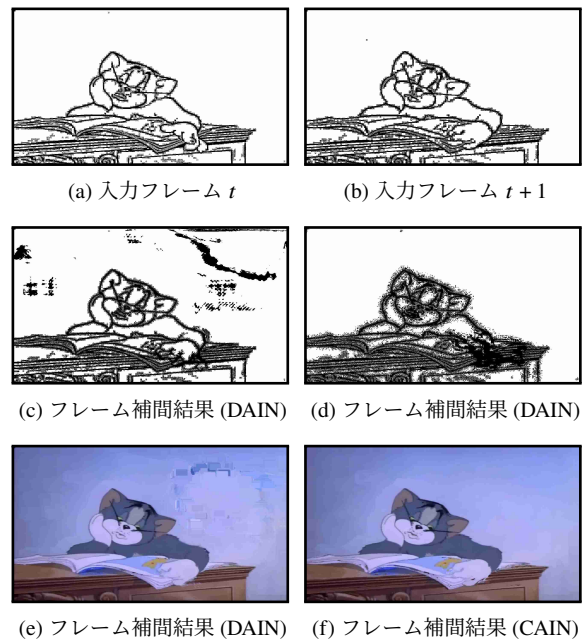


図10 フレーム補間による背景ノイズの出現例
(a) (d) は画素値 0~254 を黒画素として 2 値化