

時間相関カメラと深層学習を用いたオプティカルフロー推定に関する研究

1235069 橋本 悠衣里 【画像情報工学研究室】

Study on optical flow estimation using correlation image sensor and deep learning

1235069 HASHIMOTO Yui 【Image Processing and Informatics Lab.】

1 はじめに

オプティカルフロー (OF) とは、3次元空間での運動を画像上に射影した光学的な動きベクトルである。これまで、3次元環境推定やジェスチャー認識などに応用され積極的に研究されている分野であり、近年では深層学習を用いた OF 推定が盛んに研究されている [1][2]。また、時間相関カメラという特殊なカメラを用いて OF 推定を行っている研究もある [3]。時間相関カメラは、各画素に入射する光強度信号と参照信号との1フレーム時間相関を出力するセンサであり、例えば運動により生じる輝度変化を複素数として記録することが可能である。

そこで我々はこれに着目し、1枚の相関画像を深層学習ネットワークの入力とした OF 推定を提案した [4]。

本研究ではさらに推定精度を向上させるため、先行研究の深層学習ネットワークを改良し新たなデータセットを用いて実験を行い、学習済みモデルとの比較を行う。

2 原理

2.1 時間相関イメージセンサ

時間相関イメージセンサは、入射光と参照信号との時間相関値と強度出力を画像として出力する。入射光と参照信号との時間相関値を複素相関画像といい、参照信号として周波数 ω の複素正弦波を用いると式 (1) と表せる。

$$g(x, y) = \int_{-T/2}^{T/2} I(x, y, t) e^{-j\omega t} dt \quad (1)$$

また、強度画像は式 (1) の $\omega = 0$ で表される。

2.2 深層学習ネットワーク

CNN を用いて OF 推定を行う手法として FlowNetS が提案されている [1]。FlowNetS は、運動前後の静止画像を入力として OF 推定を行うネットワークである。本研究では、このネットワークをベースに相関画像の実部・虚部をそれぞれ異なるチャンネルとして入力するネットワークを構築した。FlowNetS では2枚のカラー画像6チャンネルを入力としているのに対し、提案法では相関画像の実部・虚部の2チャンネルを入力としている。提案するネットワークの構造を図1に示す。提案法では、FlowNetS のチャンネル数を1/4倍、1倍、2倍、3倍と増減させた4つのモデルで実験を行う。FlowNetS と提案法のモデルのチャンネル数を表1に示す。

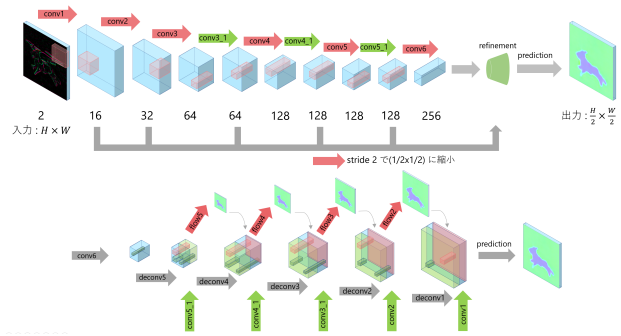


図1: 提案法ネットワークの構造
表1: 各ネットワークのチャンネル数

	FlowNetS チャンネル数		提案法							
	in	out	1/4倍		1倍		2倍		3倍	
Conv1	6	64	2	16	2	64	2	128	2	192
Conv2	64	128	16	32	64	128	128	256	192	384
Conv3	128	256	32	64	128	256	256	512	384	768
Conv3_1	256	256	64	64	256	256	512	512	768	768
Conv4	256	512	64	128	256	512	512	1024	768	1536
Conv4.1	512	512	128	128	512	512	1024	1024	1536	1536
Conv5	512	512	128	128	512	512	1024	1024	1536	1536
Conv5.1	512	512	128	128	512	512	1024	1024	1536	1536
Conv6	512	1024	128	256	512	1024	1024	2048	1536	3072
Conv6.1	1024	1024	256	256	1024	1024	2048	2048	3072	3072

FlowNetS では出力画像のサイズが入力画像の1/4サイズであるが、提案法では拡張部の層を追加し1/2サイズに変更した。さらに、学習時には式 (2) のような各層の予測フロー \hat{f}_i の誤差と重み ω_i での総和をロス関数として用い、各層のロスの計算は、式 (4) に示す予測フローと正解フロー間の全ての画素の差の2乗和である End Point Error (EPE) の総和を用いた。ここで、式 (2) 中の f, \hat{f} はそれぞれ正解フローと予測フローを表す。

$$f = \begin{bmatrix} v_x & v_y \end{bmatrix}^T, \hat{f} = \begin{bmatrix} \hat{v}_x & \hat{v}_y \end{bmatrix}^T \quad (2)$$

$$L_{EPE} = \sum_{i=1}^6 \omega_i * L_{EPE}(f, \hat{f}_i) \quad (3)$$

$$L_{EPE}(f, \hat{f}) = \sqrt{(v_x - \hat{v}_x)^2 + (v_y - \hat{v}_y)^2} \quad (4)$$

3 実験

本手法の有効性を確認するために、シミュレーションを行った。チャンネル数が1/4倍のモデルを Model1、1倍、2倍、3倍のものをそれぞれ Model2、Model3、Model4 としてそれぞれで学習を行い、評価用データセットを用いて FlyingChairDatasets により学習済みである FlowNetS(pre-traind) との比較を行った。学習の際、最適化には Adam を使用し、バッチサイズを8、学習率を0.0001と設定した。

3.1 データセット

関連画像は30枚の静止画像からマスク付き前景と後景を無作為に組み合わせ、それぞれ一様乱数により運動させ計算した。運動速度の上限は32画素/フレームとし、直線・回転運動を含んだ69,753枚の画像を生成した。また、過学習を抑えるためtrainで水平・垂直反転、0°, 90°, 180°, 270°のいずれかの角度で回転、ランダム範囲での切り取りというデータ拡張を行った。

さらに、train/testで用いる画像とは別に評価用データセットとして、運動前後のカラー画像と関連画像を新たに7830枚生成した。図2に生成した運動前後のカラー画像と関連画像を示す。



図2: 運動前後のカラー画像と関連画像

4 実験結果

testの際に最良の結果が得られた時のtrain/testの平均EPEを表2に示す

表2: 平均EPE

	train(55870枚)	test(13883枚)
Model1	0.4863	0.4015
Model2	0.2506	0.2174
Model3	0.1486	0.169
Model4	0.1839	0.1542

平均EPEは、提案法の4つのモデルの中ではModel4で最も良い精度が得られた。また、評価用データを使用して得られたOFと正解ラベルを図3に、平均EPEを表3に示す。さらに、正解ラベルとそれぞれのモデルで得られたOFとの誤差の大きさをグレイスケール画像として可視化したものを図4に示す。これは、pre-traindの最大誤差0.33で正規化を行っており、誤差0を白、誤差0.33以上を黒と表示したものになっている。

表3: 評価における平均EPE

pre-traind	Model1	Model2	Model3	Model4
0.4852	0.5042	0.3646	0.3107	0.3133

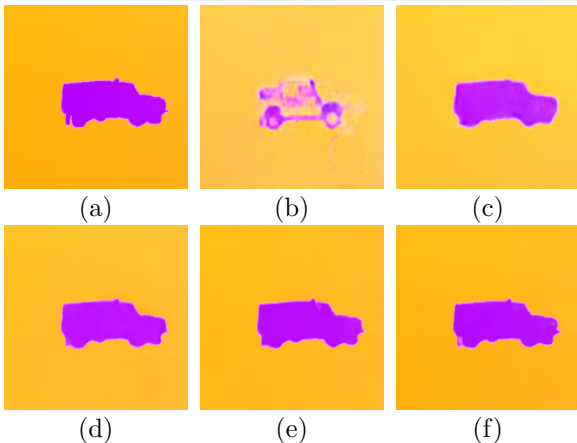


図3: OFと正解ラベル。(a)GT, (b) pre-traind, (c) Model1, (d) Model2, (e) Model3, (f) Model4

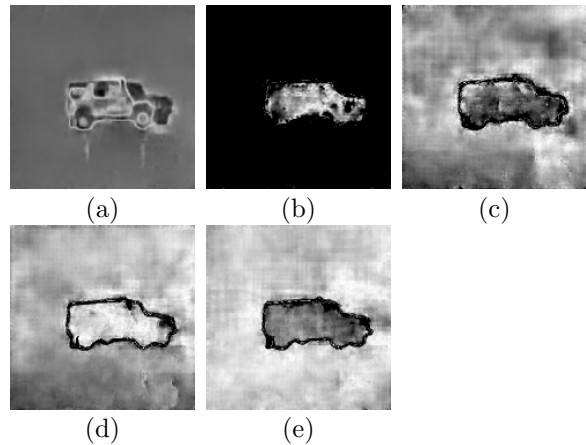


図4: 誤差の大きさをグレイスケール画像として可視化。(a)pre-traind, (b)Model1, (c)Model2, (d)Model3, (e) Model4

表2, 3より、チャンネル数をある程度増加させるにつれて推定精度が高くなっており、pre-traindと同じチャンネル数であるModel2以降では、提案法の方が高精度となっていることが確認できる。このことから、チャンネル数を増加させることで関連画像からより多くの特徴を得ることができていると言える。また図4より、pre-traindは全体的に誤差が広がっているのに対して、Model2以降の提案法は前景と後景の境界に大きな誤差が集中し、全体的に見るとpre-traindよりも誤差が少ないことが確認できる。そのため、2枚のカラー画像を1枚の関連画像に変更することで、pre-traindよりも運動に対する有益な情報が得られていると考えられる。

結果より、関連画像は運動に対する有益な情報を保持していることから、その他のネットワークも同様に入力画像を関連画像に変更することで高精度となることが期待される。

5 まとめ

本研究では1枚の関連画像を入力とした深層学習ネットワークを構築し、新たなデータセットを用いてpre-traindとの比較を行った。その結果、pre-traindと同チャンネル数であるModel2以降全てのモデルで高精度を確認することができ、関連画像の有効性を示した。

参考文献

- [1] A.Dosovitskiy, *et. al.*, "FlowNet: Learning Optical Flow with Convolutional Networks." In 2015 IEEE international Conference on Computer Vision (ICCV), pp.2758-2766, 2015.
- [2] Iig, E., *et. al.*, "FlowNet 2.0: Evolution of optical flow estimation with deep networks." In: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), pp.1647-1655, 2017.
- [3] S.Ando, *et. al.*, "Correlation image sensor: two-dimensional matched detection of amplitude modulated light." IEEE Transactions on Electron Devices, Vol.50, No.10, pp.2059-2066, 2003.
- [4] T.Kurihara, *et. al.*, "Optical flow estimation using a correlation image sensor based on FlowNet based neural network." In Proc. of 15th International Conference on Computer Vision, Theory and Applications (VISAPP), Vol.4, pp.847-852, 2020.