

## SiamMask を用いた動画修復

坪田 颯生<sup>†</sup>      岡部 誠<sup>†</sup>(正会員)      工藤 隆朗<sup>††</sup>  
 由良 俊樹<sup>†††</sup>      本間 祐作<sup>†††</sup>

<sup>†</sup>静岡大学,      <sup>††</sup>株式会社 IMAGICA GROUP アドバンストリサーチグループ,  
<sup>†††</sup>株式会社 IMAGICA Lab. 技術研究開発部

## Video Completion Using SiamMask

Satsuki TSUBOTA<sup>†</sup>, Makoto OKABE<sup>†</sup>(Member), Takaaki KUDO<sup>††</sup>,  
 Toshiaki YURA<sup>†††</sup>, Yusaku HOMMA<sup>†††</sup>

<sup>†</sup>Shizuoka University, <sup>††</sup>IMAGICA GROUP Inc., <sup>†††</sup>IMAGICA Lab. Inc.

〈あらまし〉本研究では、いずれも既存手法である SiamMask と動画修復技術を用いることで、動画内の物体やロゴ、文字、ノイズなどの消去をなるべく人手をかけずに行うシステムを開発した。ユーザは動画の最初のフレームにおいて、目的物体をバウンディングボックスで囲むことで指定する。バウンディングボックスは SiamMask への入力となる。SiamMask は目的物体を追跡しつつ、各フレームにおいて大雑把なマスクを生成する。生成されたマスクを入力として動画修復技術を用いることで修復結果が生成される。本研究が目指すのは、バウンディングボックスを描いたら即座に修復結果が生成されるようなシステムだが、現状では SiamMask によって生成される大雑把なマスクが必ずしも完璧に目的物体を捉えない場合もあり、そのような場合はユーザによる手作業での修正（フレームごとの画像編集作業）が必要となる。また本手法の応用として、入力動画と修復結果の差分に基づいて目的物体の高精度なマスクを生成する手法も提案する。

キーワード：動画修復, マスク生成, SiamMask, 物体追跡

<Summary> In this project, we are developing a system to remove objects, logos, annotations, and noises from videos with as little human intervention as possible by using SiamMask and a video completion method, both of which are existing methods. The user specifies a target object by drawing a bounding box around it in the first frame of the video; this bounding box is taken as input by SiamMask. SiamMask then tracks the target object and produces its mask in each frame. The resulting masks are then taken as input by a video completion method, which produces the final video completion result. The goal of this project is that, after drawing the bounding box, the user immediately obtains the video completion result. However, the mask produced by our current method is not always perfect. When imperfections arise, the user still has to manually modify the mask using an image editing software. As an application of this method, we also propose a method to generate a highly accurate mask of the target object based on the difference between the input video and the restoration result.

**Keywords:** video completion, mask generation, SiamMask, object tracking

## 1. はじめに

映像制作の現場において、動画内の物体やロゴ、文字、ノイズなどを消去することは、多くの需要があると同時に大変な労力を要する作業である。この作業を行うには、まず、消去したい物体を指定するためのマスクを作成する必要がある。次に、動画修復技術を適用することで、指定した物体を動画か

ら消去する。適切に作成されたマスクが与えられれば、近年開発された動画修復技術<sup>1)~3)</sup>を用いることで、物体が消去されたということがわからないような自然な動画を短時間で出力できる。しかし、物体を指定するためのマスクはすべてのフレームにおいて手作業で描く必要があり、作業に多くの時間を要する。そこで我々は、動画修復のためのマスクを容易に短時間で作成し、動画修復を効率的に行えるようなシステ

ムを開発している。さらに、元の動画と修復結果の差分に基づき、物体の高精度なマスクを生成する実験も行ったので報告する。

画像や動画から物体を正確にセグメンテーションするため数多くの手法が提案されてきた。その多くはユーザの入力無しに画像から物体を自動的に抽出する技術<sup>4)</sup>や、それらを動画へ拡張した技術<sup>5)</sup>などであるが、我々はあらかじめ学習済みの物体ではなく、ユーザによってその場で指定された物体のマスクを作成することを想定しているのでこれらの技術を適用することはできない。そこで我々は使い勝手が良く、計算効率も良い SiamMask<sup>6)</sup>に注目した。ユーザが動画の最初のフレームで物体をバウンディングボックスで囲むと、SiamMask はそれ以降のフレームでその物体を追跡しつつ、各フレームにおいてマスクを作成する。

我々の目的は、ユーザが物体をバウンディングボックスで囲んだら、その物体が消去された動画を即座に出力することである。つまり、我々の目的は物体のマスクを正確に作成することではなく、動画修復をより正確により効率良く行うことである。動画修復に用いられるマスクは目的物体の形に完璧に一致する必要はないが、目的物体を完全に覆うような形でなければならない。そこで我々は、SiamMask で得られるマスクを膨張させて一回り大きくなるようにする。しかし、マスクを膨張させても常に完璧に目的物体を覆えるわけではないため、その場合はユーザが手作業で追加のマスクを描いて修正する必要がある。

## 2. 提案手法

例として、図1に提案手法を用いて動画からローラーブレードに乗っている人物を消去する際の作業工程の概要を示す。まず、ユーザは動画の最初のフレームで人物を囲むようなバウンディングボックスを描く(図1-Inputの赤枠)。SiamMaskはこのバウンディングボックスを入力として、各フレームにおけるマスクを生成する(図1-SMの赤いピクセル)。このマスクの生成の処理には数秒を要する。SiamMaskにより生成されるマスクは目的物体の大部分を覆うが、覆い切れない部分も多く、すなわち、目的物体がマスクをはみ出してしまふことがある。例えば、6フレーム目では手や足など、人物の領域の端に当たる部分を覆い切れていない(図1-SMの黄矢印)。我々はそれらの部分をマスクが覆うようにするために、各フレームのマスクを一回り大きくするような膨張操作(モルフォロジー変換)を加えた(図1-Dilated)。しかしそれでも目的物体を覆い切れていない箇所がいくつかある。例えば、10フレーム目や33フレーム目では手が覆われていない(図1-Dilatedの緑矢印)。このような場合、ユーザはGIMP<sup>7)</sup>などの画像編集ツールを使って手作業でマスクを修正する必要がある(図1-Modified)。最終的に修正されたマスクは動画修復技術<sup>3)</sup>へ入力され、動画からローラーブレードに乗っている人物を消去することができる(図1-Result)。この動画修復

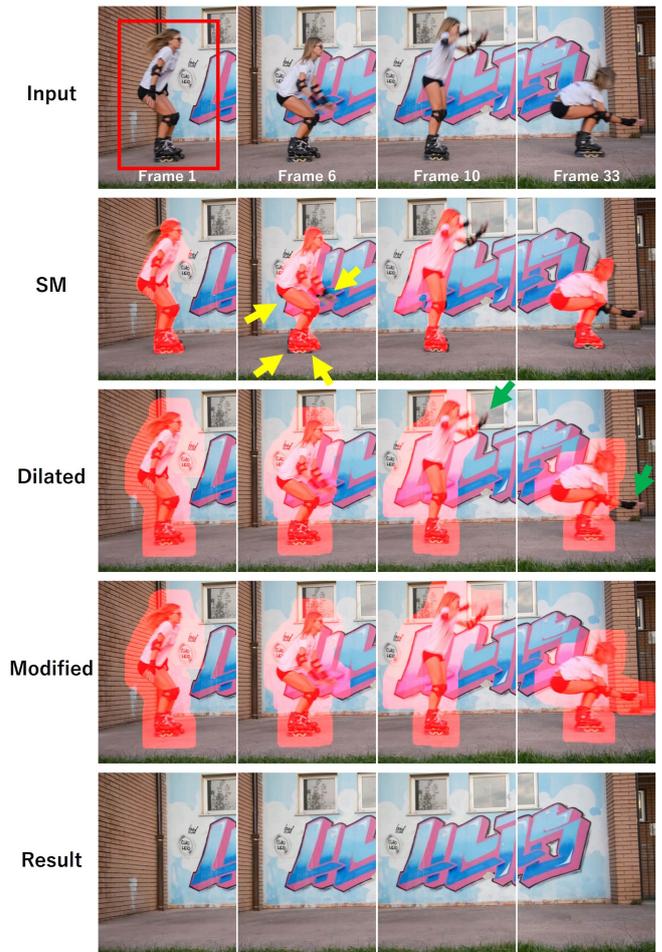


図1 提案手法の概要

Fig. 1 The overview of the proposed method

の処理には数秒を要する。

## 3. 結果と考察

3種類の動画を対象として提案手法を適用し、生成した動画について12名の学生を被験者として実験を行った。

### 3.1 実験結果

図2に提案手法による動画修復結果と背景差分によって生成されたマスクを示す。なお、マスク生成については3.2節で後述する。

“soccerball”の動画はフレーム数が48である。その12フレーム目を図2の1行目に示す。この動画については、本研究の目的を達成することができた。すなわち、ユーザが動画の最初のフレームにおいてバウンディングボックスでボールを囲めば、即座にボールが消去された動画を得ることができた。

“bmx-trees”の動画はフレーム数が80である。その31フレーム目を図2の2行目に示す。動画の最初のフレームでバウンディングボックスを描くと、SiamMaskと膨張操作によって72枚のフレームで適切なマスクを得たが、残り8枚のフレームでは正確なマスクが得られなかった。それらのフレームのマスクをGIMPを用いて手作業で修正することで満足のいく修復結果を得ることができた。



図2 提案手法による動画修復結果と背景差分によって生成されたマスク  
 Fig. 2 Video completion results by the proposed method and the mask generated by the background subtraction

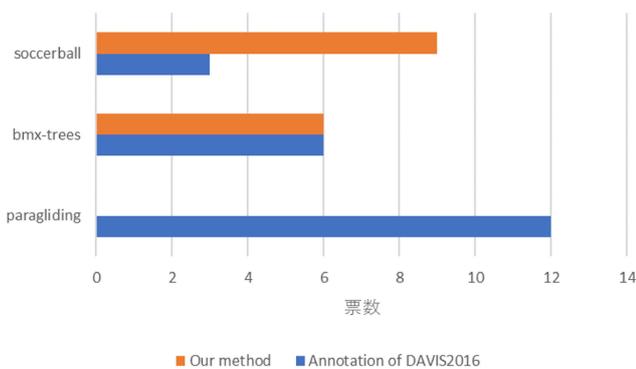


図3 定性的な評価の結果  
 Fig. 3 Results of the qualitative evaluation

“paragliding”の動画はフレーム数が70である。その31フレーム目を図2の3行目に示す。この動画ではSiamMaskと膨張操作によって61枚のフレームで適切なマスクを得て、残り9枚のフレームではGIMPでの手作業の修正により適切なマスクを得た。

生成した動画に対し、定性的な評価を行った。被験者は12名の学生であり、比較対象としてDAVIS2016<sup>8)</sup>のアノテーションであるマスクに膨張操作をして動画修復技術により動画修復を行ったものを用意した。各動画について入力動画と動画修復結果を被験者に見せ、提案手法と比較対象のうち自然に物体を消去できていると思った方に投票してもらい、投票した理由も答えてもらった。

結果は図3のようになった。ここでオレンジの棒グラフは提案手法を選択した人数、青の棒グラフは既存手法を選択した人数を示す。“soccerball”と“bmx-trees”に関しては投票結果にばらつきがあり、被験者の大半の意見として本手法と比較対象の結果にほとんど違いは見られないとのことであった。比較対象が全フレームにおいて手作業で作られた正確なマスクであるのに対し、本手法では“soccerball”につい

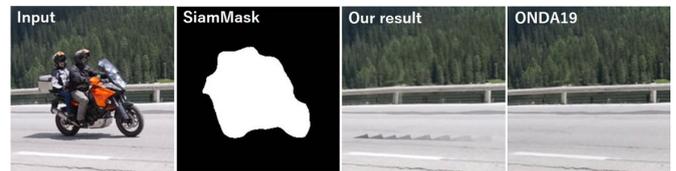


図4 提案手法では影の消去ができない  
 Fig. 4 Shadow elimination is not possible with the proposed method

ては手作業無し、“bmx-trees”については8枚のフレームでの手作業により同程度のクオリティの修復動画を生成したことになる。これは動画修復の作業において、提案手法を用いることで従来よりも大幅な負担軽減になっていると言える。“paragliding”では被験者全員がDAVIS2016のアノテーションによる修復動画に投票し、その理由として本手法による修復動画には物体の一部が少し残っているということを挙げた。“paragliding”の動画にはパラシュート部分と人間の間にロープがある。DAVIS2016のアノテーションではこのロープ部分もマスクされているが、一方、SiamMaskはこのロープのような細い物体のマスクを作る能力がない。そのためロープが消去されず動画に残ってしまった。

今後の目標は、SiamMaskを改良することで、さらに適切なマスクを生成できるようなシステムを開発し、ユーザの手作業での修正といった負担を減らしていくことである。今回提案した手法では、図1の手や、図2の3行目のパラグライダーのような小さい物体や細い物体のマスクを適切に生成することは難しいことがわかった。また、図2の2行目の木の後ろから急に現れる自転車のように、動画内で突然消えたり出現したりする物体を捉えてマスクを生成することも困難であることがわかった。これは、SiamMaskがマスクの推定をする際に、直前のフレームで生成されたマスクの情報を使うためだと考えられる。これらの問題を今後解決していきたいと思う。

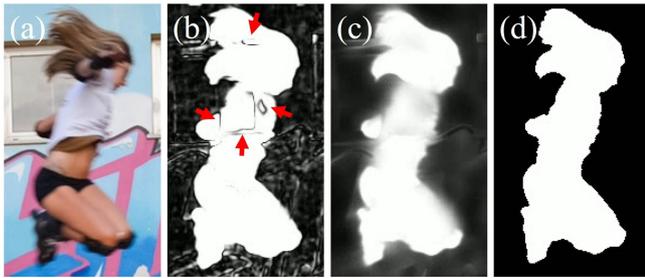


図5 背景差分によるマスクの生成プロセス

Fig. 5 The process of generating a mask by background subtraction

さらに別の問題として、提案手法では物体の下にできる影を取り除くことも難しい。例として図4を示す。目的の物体はオートバイで、その下にくっきりとした影があるが(図4-Input), SiamMaskでは正確に捉えることができない(図4-SiamMask)。これはSiamMaskが、画像分類のためのデータセットを事前学習したネットワークをバックボーンに使用しているからではないかと推察する。つまり、画像分類のためには影を考慮する必要がない。影の部分をほとんど捉えていないため、膨張操作を行っても影を消すことはできなかった(図4-Our result)。影を捉え、既存手法<sup>3)</sup>で得られたような結果(図4-ONDA19)を達成することも今後の目標の1つである。

### 3.2 背景差分によるマスク生成への応用

動画修復に使うマスクは上述のように目的物体の形に完璧に一致するマスクである必要はないが、そのようなマスクを作成することは、映像制作の現場における別の重要な問題の1つである。そこで我々は、元の動画と修復結果の差分に基づき、SiamMaskでは得られなかったような高精度なマスクを生成できるのではないかと考えた。これが上手くいけば、ユーザは動画修復に必要な程度のラフなマスクを作成するだけで高精度なマスクを手に入れることができる。

背景差分によりSiamMaskより高精度なマスクを生成するプロセスを図5に示す。まず、図1-Inputと図1-Resultの差分によって図5(b)のような画像が得られる。人物の領域にマスクができてはいるが、ところどころに黒い穴が開いていて(図5(b)の赤矢印)、このままでは高精度なマスクと言えない。図5(b)を観察すると、ここに見られる穴は主に背景、つまり、図1-Resultに見られる窓枠や壁の模様が原因であることがわかる。つまり、背景差分におけるノイズは背景画像のエッジに起因する傾向が強い。そこで、cross bilateral filter<sup>9)</sup>を適用し、入力動画(図5(a))に見られるエッジは保存し、それ以外のエッジを消去するように平滑化することで、図5(c)を得る。最後に閾値によって2値化した画像が図5(d)である。以上のプロセスにより、図6のようなマスクを得ることができる。

図2の5列目は、この手法により生成したマスクである。これは、図2の2列目のSiamMaskによるマスクよりも高精



図6 背景差分によるマスクの生成結果

Fig. 6 Masks generated by background subtraction

度なマスクであることがわかる。動画修復に必要な程度のラフなマスクは高精度なマスクに比べて圧倒的に簡単に作成できるので、本手法によって高精度なマスクを必要とするユーザの負担をかなり軽減できているのではないかと期待している。将来はユーザテストや現場で働くデジタルアーティストへのインタビューを実施することによって実際の効果を確認したい。

## 4. おわりに

本稿ではいずれも既存手法であるSiamMaskと動画修復技術を用いることで、素早く簡単に動画修復を行えるようなシステムを提案した。実験の結果、提案手法を用いるとバウンディングボックスを描くだけで動画から物体を完全に消去できる場合があることがわかった。また手作業での修正が必要な場合も、編集すべきフレーム数を大幅に削減できることがわかった。さらにユーザ調査の結果、提案手法の動画修復結果は、正確に作られたDAVIS2016<sup>8)</sup>のマスクによる結果と比較しても、遜色なくオリティであることが確認できた。

しかし、ロープのように細い物体や動画内に突然現れる物体、そして物体の影などに対してはSiamMaskでは適切にマスクを作ることができず高精度な動画修復を行うことができなかった。今後はこれらの問題を解決していきたいと思う。

## 参考文献

- 1) A. Bokov, D. Vatolin: "100+ Times Faster Video Completion by Optical-Flow-Guided Variational Refinement", Proc. of IEEE International Conference on Image Processing 2018, pp. 2122-2126 (2018).
- 2) R. Murase, Y. Zhang, T. Okatani: "Video-Rate Video Inpainting", Proc. of IEEE Workshop on Applications of Computer Vision 2019, pp. 1553-1561 (2019).
- 3) M. Okabe, K. Noda, Y. Dobashi, K. Anjyo: "Interactive Video Completion", IEEE Computer Graphics and Applications (CG&A), Vol. 40, pp. 127-139 (2019).
- 4) R. Girshick, I. Radosavovic, G. Gkioxari, P. Dollár, K. He, Detectron, <https://github.com/facebookresearch/detectron> (2018).
- 5) P. Voigtlaender, M. Krause, A. Ošep, J. Luiten, B. B. G. Sekar, A. Geiger, B. Leibe: "MOTS: Multi-Object Tracking and Segmentation", Computer Vision and Pattern Recognition (2019).
- 6) Q. Wang, L. Zhang, L. Bertinetto, W. Hu, P. H. Torr: "Fast

Online Object Tracking and Segmentation: A Unifying Approach”, Proc. of Computer Vision and Pattern Recognition 2019 (2019).

- 7) The GIMP Development Team, Gimp, <https://www.gimp.org>
- 8) F. Perazzi, J. Pont-Tuset, B. McWilliams, L. V. Gool, M. Gross, A. Sorkine-Hornung: “A benchmark Dataset and Evaluation Methodology for Video Object Segmentation”, Computer Vision and Pattern Recognition 2017 (2017).
- 9) S. Paris, F. Durand: “A Fast Approximation of the Bilateral Filter Using a Signal Processing Approach”, Proc. of European Conference on Computer Vision 2006, pp. 568-580 (2006).

(2020年12月17日 受付)

(2021年4月17日 再受付)



### 坪田 颯生

2020年 静岡大学工学部数理システム工学科卒業。現在、静岡大学大学院総合科学技術研究科修士課程に在学。動画中の人や物体のマスク生成に関する研究に取り組んでいる。



### 岡部 誠 (正会員)

2008年 東京大学大学院情報理工学系研究科博士課程修了。博士(情報理工学)。2008年マックスプランク研究所ポストドクター。2010年電気通信大学助教。2016年静岡大学助教。2017年同准教授。動画データの分析と映像製作ツールに関する研究に取り組んでいる。



### 工藤 隆朗

2005年 東京農工大学工学府電気電子工学専攻修了(修士)。同年、株式会社IMAGICAに入社。主に映像圧縮技術の改善、パッケージメディア生産ラインの構築に従事。2021年より株式会社フォトロンにて先端映像技術の研究に取り組んでいる。



### 由良 俊樹

2006年 上智大学理工学研究科修士課程修了。同年、株式会社IMAGICAに入社。主に映像作品の色管理や画像処理技術の研究に従事。2021年より株式会社フォトロンにて先端映像技術の研究に取り組んでいる。



### 本間 祐作

2018年 首都大学東京大学院システムデザイン研究科経営システムデザイン学域修士課程修了。同年、株式会社IMAGICAに入社。2021年より株式会社フォトロンにてクラウド関連の事業や先端映像技術の研究に取り組んでいる。