



**HAL**  
open science

## Object removal from complex videos using a few annotations

Thuc Trinh Le, Andrés Almansa, Yann Gousseau, Simon Masnou

### ► To cite this version:

Thuc Trinh Le, Andrés Almansa, Yann Gousseau, Simon Masnou. Object removal from complex videos using a few annotations. *Computational Visual Media*, 2019, 5, pp.267-291. 10.1007/s41095-019-0145-0 . hal-02168653

**HAL Id: hal-02168653**

**<https://hal.science/hal-02168653v1>**

Submitted on 28 Jun 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Object removal in complex videos from a few annotations

Thuc Trinh Le<sup>1</sup> (✉), Andrés Almansa<sup>2</sup>, Yann Gousseau<sup>1</sup>, and Simon Masnou<sup>3</sup>

© The Author(s) 2015. This article is published with open access at Springerlink.com

**Abstract** We present a system for the removal of objects from videos. As an input, the system only needs a user to draw a few strokes on the first frame, roughly delimiting the objects to be removed. To the best of our knowledge, this is the first system allowing the semi-automatic removal of objects in videos with complex backgrounds. The key steps of our system are the following: after initialization, segmentation masks are first refined and then automatically propagated through the video. The missing regions are then synthesized using video inpainting techniques. Our system can deal with multiple, possibly crossing objects, with complex motions, and with dynamic textures. This results in a computational tool that can alleviate tedious manual operations for editing high-quality videos.

**Keywords** objects removal, objects segmentation, object tracking, video inpainting, video completion.

## 1 Introduction

In this paper, we propose a system to remove one or several objects from a video, starting with only a few user annotations. More precisely, the user only needs to approximately delimit in the first frame the objects to be edited. Then, these annotations are refined and propagated through the video. One or several objects can then be removed automatically.

This results in a flexible computational video editing tool, with numerous potential applications. Removing unwanted objects (such as a boom microphone) or people (such as an unwanted wanderer) is a common task in video post-production. Such tasks are critical given the time constraints of movie production and the prohibitive costs of reshooting complex scenes. They are usually achieved through extremely tedious and time-consuming frame-by-frame processes, for instance using the Rotobrush tool from Adobe After Effects [2] or professional visual effects softwares such as SilhouetteFX or Mocha. More generally, the proposed system paves the way to sophisticated movie editing tasks, ranging from crowd suppression to unphysical scenes modifications, and has potential applications for multi-layered video editing.

Two main challenges arise in developing such a system. First, not a single part of the objects to be edited shall be left over in the tracking part of the algorithm; otherwise, they are propagated and enlarged by the completion step, resulting in unpleasant artifacts. Second, our visual system is good at spotting temporal discontinuities and aberrations, making the completion step a tough one. We address both these issues in this work.

The first step of our system consists of transforming a rough user annotation into a mask that accurately represents the object to be edited. For this, we use a classical strategy relying on a CNN-based edge detector, followed by a watershed transform yielding super-pixels, which are eventually selected by the user to refine the segmentation mask. After this step, a label is then given to each object. The second step is the temporal propagation of the labels. There we make use of state-of-the-art advances in CNN-based multiple objects segmentation. Besides, our approach includes an original and crucial algorithmic brick which consists in learning the transition zones between objects and the background, in such a way that the objects will be fully covered by the propagated masks. We call

1 LTCI, Télécom ParisTech, Université Paris-Saclay, 75013 Paris, France . E-mail: thuc.le@telecom-paristech.fr, yann.gousseau@telecom-paristech.fr.

2 MAP5, CNRS & Université Paris Descartes, 75006 Paris, France. E-mail: andres.almansa@parisdescartes.fr.

3 Univ Lyon, Université Claude Bernard Lyon 1, CNRS UMR 5208, Institut Camille Jordan, 69622 Villeurbanne, France. E-mail: masnou@math.univ-lyon1.fr.

Manuscript received: 2019-03-15;

the resulting brick a *smart dilation* by analogy with the dilation operators of mathematical morphology. Our last step is then to remove some or all of the objects from the video, depending on the user's choice. For this, we employ two strategies: a motion-based pixel propagation for static background and a patch-based video completion for dynamic textures. Both methods rely heavily on the knowledge of segmented objects. This interplay between objects segmentation and the completion scheme improves the method in many ways: it allows for better video stabilization, for a faster and more accurate search for similar patches, and for a more accurate foreground/background separation. These improvements yield completion results with very little or no temporal incoherence.

We illustrate the effectiveness of our system through several challenging cases including severe camera shake, complex and fast object motions, crossing objects, and dynamic textures. We evaluate our method on various datasets, in both objects segmentation and objects removal. Moreover, we show on several examples that our system yields comparable or better results than state-of-the-art video completion methods applied on manually segmented masks. This paper is organized as follows: First, we briefly explore some related works (section 2). Next, we introduce our proposed approach which includes three steps: First frame annotation, objects segmentation and objects removal (section 3). Finally, we show experimental results as well as some evaluation and comparison with other state-of-the-art methods. A shorter version of this work can be found in [40].

## 2 Related works

The proposed computational editing approach is related to several families of works that we now briefly review.

### 2.1 Video object segmentation

Video object segmentation, the process of extracting space-time segments corresponding to objects, is a widely studied topic whose complete review is beyond the scope of this paper. For a long time, such methods have not been accurate enough to avoid using green-screen compositing to extract objects from videos. Significant progress for the supervised segmentation has been achieved by the end of the 2000s, see e.g. [2], and in particular, the use of supervoxels became the most flexible way to incorporate user annotations in the segmentation process [44, 78]. Other efficient approaches to the supervised object segmentation

problem are introduced in [49, 53].

A real breakthrough occurred with approaches relying on Convolutional Neural Networks (CNN). In the DAVIS-2016 challenge [63], the most efficient methods were all CNN-based, both for the unsupervised and semi-supervised tasks. For the semi-supervised task, where a first frame annotation is available, methods mostly differ in the way they train the networks. The One Shot Video Object Segmentation (OSVOS) method, introduced in [8], starts from a pre-trained network and retrains it using a large video dataset, before fine-tuning it per-video using the annotation at the first frame to focus on the object being segmented. With a similar approach, [62] relies on an additional mask layer to guide the network. The method in [7] further improves the results from OSVOS with the help of Multi Networks Cascade (MNC) [20].

All these approaches work image-per-image without explicitly checking for temporal coherence, and therefore can deal with large displacements and occlusions. However, since their backbone is a network used for semantic segmentation, they cannot distinguish between instances of the same class or between objects that resemble each other.

Another family of works deals with the segmentation of multiple objects. Compared with the single object segmentation problem, an additional difficulty here is to distinguish between different object instances which may have similar colors and may cross each other. Classical approaches include graph-based segmentation using color or motion information [42, 58, 87], the tracking of segmentation proposals [17, 45], or bounding box guided segmentation [22, 71].

The DAVIS 2017 challenge [66] established a ranking between methods aiming at the semi-supervised segmentation of multiple objects. Again, the most efficient methods were CNN-based. It is proposed in [77] to modify the OSVOS network [8] to work with multiple labels and to perform online fine-tuning to boost the performances. In [37], the networks introduced in [74] are adapted to the purpose of multiple objects segmentation through the heavy use of data augmentation, still using annotation of the first frame. The authors of this work also exploit motion information by adding optical flow information to the network. This method is further improved in [46] by using a deeper architecture and a re-identification module to avoid propagating errors. This last method has achieved the best performance in the DAVIS-2017 challenge [66]. With a different approach, Hu et al. [31]

employ a recurrent network exploiting the long-term temporal information.

Recently, with the release of a large-scale video object segmentation dataset called YouTube Video Object Segmentation (YouTube-VOS) [83], many further improvements have been made in the field. Among them, one of the most notable work is PreMVOS [48] which has won the 2018 DAVIS Challenge [9] and Youtube-VOS challenge [83].

In PreMVOS, the algorithm first generates a set of accurate segmentation mask proposals for all objects in each frame of a video. To achieve this, a variant of the Mask R-CNN [29] object detector is used to generate coarse object proposals, then a fully convolutional refinement network inspired by [82] and based on the DeepLabv3+ [11] architecture produces accurate pixel masks for each proposal. Secondly, these proposals are selected and merged into accurate and temporally consistent pixel-wise object tracks over the video sequence. In contrast with PreMVOS which focuses on the accuracy, some methods trade off accuracy for speed. Those methods take the first frame with its mask annotation either as guidance to slightly adjust parameters of the segmentation model [85] or as a reference for segmenting the following frames without tuning the segmentation model [12, 13, 57].

Although these methods yield impressive results in terms of the accuracy of the segmentation, they may not be the optimal solutions for the problem we consider in this paper. As said above, when removing objects from a video it is crucial for the video completion step that no part of the removed objects remains after the segmentation. Said differently, we are in a context where *recall* is much more important than *precision*, see Section 4.2 for the definitions of these metrics. In the experimental section, we compare our segmentation approach to several state-of-the-art methods with the aim of optimizing a criterion which penalizes under-detection of objects.

## 2.2 Video editing

Recently, advances in both the analysis and the processing of videos have permitted advances in the emerging field of computational video editing. Examples include, among others, tools for the automatic, dialogue-driven selection of scenes [41], time slice video synthesis [19], or methods for the separate editing of reflectance and illumination components [6]. It is proposed in [89] to identify accurately the background in videos to either improve the stabilization process or proceed to tasks such as background

suppression or multi-layered editing. In a sense, our work is more challenging since we need to identify moving objects with enough accuracy so they can be removed seamlessly.

Because we learn a transition zone between objects and the background, our work is also related to image matting techniques [43], and their extension to videos [16] as a necessary first step for editing and compositing tasks. Lastly, since we deal with semantic segmentation and multiple objects, our work is also related to the soft semantic segmentation recently introduced for still images [1].

## 2.3 Video inpainting

Image inpainting, also called image completion, refers to the task of reconstructing missing or damaged image regions by taking advantage of the image contents outside these missing regions.

The first approaches were variational [50], or PDE-based [4] and dedicated to the preservation of geometry. They were followed by patch-based methods [18, 23], inherited from texture synthesis methods [24]. Some of these methods have been adapted to videos, often by mixing pixel-based approaches for reconstructing the background and greedy patch-based strategies for moving objects [60, 61]. In the same vein, different methods have been proposed to improve or speed up the reconstruction of the background [26, 30], with the strong limitation that the background should be static. Other methods yield excellent results in restricted cases, such as the reconstruction of cyclic motions [36].

Another family of works which performs very well when the background is static relies on motion-based pixel propagation. The idea is to first infer a motion field outside and inside the missing regions. Using the completed motion field, pixel values from outside the missing region are then propagated inside it. For example, Grossauer *et. al* describes in [28] a method for removing blotches and scratches in old movies using optical flow. A limitation of this work is that the estimation of the optical flow suffers from the presence of the scratches. Using a similar idea, but avoiding calculating the optical flow directly in the missing regions, several methods try to restore the motion field inside these missing regions by gradually propagating motion vectors [51], by sampling spatial-temporal motion patches [72, 73], or by interpolating the missing motion [5, 88].

In parallel, it was proposed in [79] to address the video inpainting problem as a global patch-based optimization problem, yielding unprecedented time

coherence at the expense of very heavy computational costs. The method in [54] was developed from this seminal contribution, by accelerating the process and taking care of dynamic texture reconstruction. Other state-of-the-art strategies rely on a global optimization procedure, taking advantage of either shift-maps [27] or an explicit flow field [32]. This last method arguably has the best results in terms of temporal coherence, but since it relies on two-dimensional patches, it is not suitable for the reconstruction of dynamic backgrounds. Recently, it was proposed in [39] to improve the global strategy of [54] by incorporating the optical flow in a systematic way. This approach has the ability to reconstruct complex motions as well as dynamic textures.

Let us add that the most recent approaches to image inpainting rely on convolutional neural networks and have the ability to infer elements that are not present in the image at hand [33, 59, 76]. To the best of our knowledge, such approaches have not been adapted to videos because their training cost is prohibitive.

In this work, we will propose two complementary ways to perform the inpainting step needed to remove objects in videos. A first method is fast and relies on a frame-by-frame completion of the optical flow, followed by the propagation of voxel values. This approach is inspired by the recently introduced method [5], itself sharing ideas with the approach from [32] and yielding impressive gains in terms of computational times. Such approaches are computationally efficient but not able to deal with moving backgrounds and dynamic textures. For these complex cases, we rely on a more sophisticated (and much slower) second approach extending the ideas we initially developed in [39].

### 3 Proposed method

The general steps of our method are as follows:

- (a) First, the user draws a rough outline of each object of interest in one or several frames, for instance in the first one (Section 3.1);
- (b) These approximate outlines are refined by the system, then propagated to all remaining frames using different labels for different objects (Section 3.2);
- (c) If some errors are detected, the user may manually correct them in one or several frames (using step (a)) and propagate these edits to the other frames (using step (b));

- (d) Finally, the user selects which of the selected objects he/she wants to remove, and the system removes the corresponding regions in the whole video, reconstructing the missing parts in a plausible way (Sections 3.3.1 and 3.3.2). For this last step two options are available : a fast one for static background and a more involved one for dynamic backgrounds.

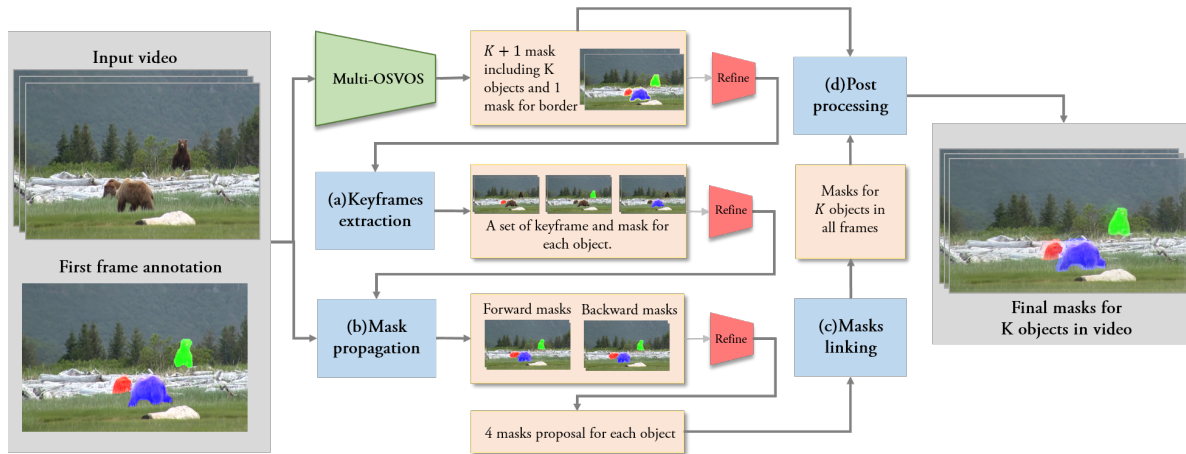
In the first step most methods only select the object to be removed. There are, however, several advantages to tracking multiple objects with different labels:

1. It gives more freedom to the user for the inpainting step with the possibility to produce various results depending on which objects are removed; in addition, objects which are labeled but not removed are considered as important by the system and therefore better preserved during the inpainting of other objects.
2. It may produce better segmentation results than tracking a single object, in particular when several objects have similar appearance.
3. It facilitates video stabilization and therefore increases the temporal coherence during the inpainting step, as shown in the results (Section 4.3).
4. It is of interest for other applications, e.g., action recognition or scene analysis.

The illustration of these steps can be found in the supplementary website <https://object-removal.telecom-paristech.fr/>

#### 3.1 First frame annotation

A classical method to cut out an object in a frame involves commercial tools such as the Magic Wand of Adobe Photoshop which is fast and convenient. However, this classical method requires many refinement steps and is not accurate with complex objects. To increase the precision and reduce the user's intervention, many methods have been proposed where interactive image segmentation is performed using scribbles, point clicks, superpixels, etc. Among them, some state-of-the-arts annotators achieve a high degree of precision by using edge detectors to find the contour map and create a set of object proposals from this map [35]; the appropriate regions are then selected by the user using point clicks. The main drawbacks of these approaches are a large computation time and a weak level of user input.



**Fig. 1** General pipeline of our object segmentation method. Given the input video and annotations in the first frame, our algorithm alternates two CNN-based semantic segmentation steps (multi-OSVOS network in green and Refining network in red) with 4 video-tracking steps (depicted as blue blocks): (a) keyframe extraction, (b) mask propagation, (c) mask linking and (d) post processing. These steps are detailed in Section 3.2.

In order to balance between human effort and accuracy, we adopt a fast and simple algorithm. Our system first generates a set of superpixels from the first image, then the user can select suitable superpixels by simply drawing a coarse contour around each object. The set of superpixels is created using an edge-based approach. More precisely, an FCN-based edge detector network introduced in [80] is applied to the first image, and its output is a probability map of edges. Superpixels are extracted from this map by the well-known watershed transform [52], which runs directly on edge scores. There are two main advantages of using this CNN-based method to compute the edge map:

1. It has shown superior performances over traditional boundary detection methods that use local features such as colors and depths. In particular, it is much more accurate.
2. It is extremely fast: one forward pass of the network takes about 2 ms hence the annotation step is performed in real time and very interactively.

After computing all superpixels, the user selects the suitable ones by drawing a contour around each target object to get rough masks. Superpixels which overlap these masks by more than 80 percent are selected. The user can also refine the mask by adding or removing superpixels using mouse clicks. As a result, accurate masks for all objects of interest are extracted in a frame within few seconds of interactive annotation.

### 3.2 Objects segmentation

In this step, we start from the object masks computed on the first frame using the method described in the previous section, and we aim at inferring a full space-time segmentation of each object of interest in the whole video. We want our segmentation to be as accurate as possible, in particular without false negatives.

Doing this in complex videos with several objects which occlude each other is an extremely challenging task. As described in Section 2, CNNs have made important breakthroughs in semantic image segmentation with extensions to video segmentation in the last two years [9, 64, 66]. However, current CNN-based semantic segmentation algorithms are still essentially image-based, and do not take global motion information sufficiently into account. As a consequence, semantic segmentation algorithms cannot deal with sequences where: (a) several instances of similar objects need to be distinguished; and (b) these objects may eventually cross each other. Examples of such sequences are *Les Loulous*<sup>1</sup> introduced in [54] or *Museum* and *Granados-S3*<sup>2</sup> introduced in [26, 27].

On the other hand, more classical video tracking techniques like optical-flow based propagation or global graph-based optimization do take global motion information into account [84]. Nevertheless, they are most often based on bounding boxes or rough descriptors and do not provide a precise delineation of objects' contours. Two recent attempts to adapt video-

<sup>1</sup>[https://perso.telecom-paristech.fr/gousseau/video\\_inpainting/](https://perso.telecom-paristech.fr/gousseau/video_inpainting/)

<sup>2</sup><http://gvv.mpi-inf.mpg.de/projects/vidinp/>

tracking concepts to provide a precise multi-object segmentation [68, 75] fail completely when objects cross each other like in the *Museum*, *Granados-S3* or *Loulous* sequences.

In the rest of this section, we describe a novel hybrid technique which combines the benefits of classical video tracking with those of CNN-based semantic segmentation. The structure of our hybrid technique is shown in Figure 1. CNN-based modules are depicted in green and red, and their inner structure is described in Section 3.2.1 and Figure 2. Modules that are inspired from video-tracking concepts are depicted in blue and are detailed in Section 3.2.2.

Note that the central part of Figure 1 operates in a frame-by-frame basis. Each segmentation proposal by the *Multi-OSVOS network* (in green), or by the *Mask propagation* module (in blue) is improved by the *Refinement network* (in red). In the right part of the figure the *Mask linking* module (in blue) builds a graph that links all segmentation proposals from the previous steps, and makes a global decision on the optimal segmentation for each of the  $K$  objects to be tracked. Finally the *Keyframe extraction* module is required to set sensible temporal limits to the *Mask propagation* iterations, and the final *post-processing* module further refines the result with the objective of maximizing the recall, which is much more important than precision in the case of video inpainting. All these modules will be explained in more detail in the next sections.

### 3.2.1 Semantic segmentation networks

Our system uses two different semantic segmentation networks: a *multi-OSVOS* network and a *refinement* network. Both operate on a frame by frame basis.

Our implementation of *multi-OSVOS* computes  $K+1$  masks for each frame:  $K$  masks for the  $K$  objects of interest and one novel additional mask covering the objects' boundaries. We call this latter mask a *smart dilation* layer, it is a key to guarantee that the segmentation does not miss any part of the objects, which is especially difficult in the presence of motion blur.

While the *multi-OSVOS network* provides a first prediction, the *refinement network* takes mask predictions as an additional guidance input and improves those predictions based on image content, similarly to [62].

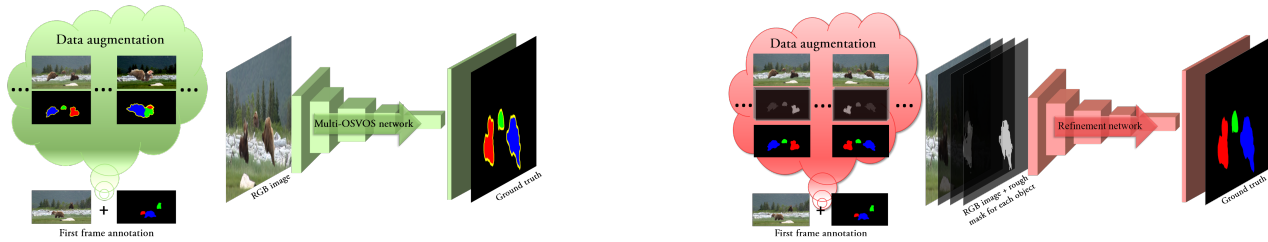
Training these networks is a challenging task, because the only labeled example we can rely on (for supervised training) is the first annotated frame and the corresponding  $K$  masks. The next paragraphs focus

on our networks' architectures and on semi-supervised training techniques that we use to circumvent the training difficulty.

**Multi-OSVOS network.** The training technique of our semantic segmentation networks is mainly inspired from the OSVOS network [8], a breakthrough which achieved the best performance in DAVIS-2016 challenge [63]. The OSVOS network uses a transfer learning technique for image segmentation: the network is first pre-trained on a large database of labeled images. After training, this so-called *parent* network can roughly separate all foreground objects from the background. Next, the parent network is fine-tuned using the first frame annotation (annotation mask and image) in order to improve the segmentation of a particular object of interest. OSVOS has proven to be a very fast and accurate semi-supervised method to obtain a background/foreground separation. Our Multi-OSVOS network uses a similar transfer learning technique, yet with several important differences:

- Our network can identify different objects separately (instead of a simple foreground/background segmentation) and provides a smart dilation mask, i.e. a smart border which covers the interfaces between segmented objects and the background, and reduces a lot the number of false negative pixels. The ground truth for this smart dilation mask is defined in the fine-tuning step by a 7-pixels wide dilation of the union of all object masks.
- Unlike OSVOS, which uses a fully convolutional network (FCN) [47], our network uses the Deeplab v2 [10] architecture as the parent model since it outperforms FCN in some common datasets such as PASCAL VOC 2012 [25].
- In the fine-tuning training step we adopt a data augmentation technique in the spirit of Lucid Tracker [37]: we remove all objects from the first frame using Newson et al's image inpainting algorithm [55], then the removed objects undergo random geometric deformations (affine and thin plate deformations), and eventually they are Poisson blended [65] over the reconstructed background. This is a sensible way of generating large amounts of labeled training data with an appearance similar to what the network might observe in the following frames.

The smart dilation mask is of particular importance to ensure that segmentation masks do not miss any part of the object, which is typically difficult in the presence of



**Fig. 2** Two networks used in the general pipeline presented in Figure 1. Left: multi-OSVOS network, Right: refinement network. They serve different purposes: the multi-OSVOS network helps us separating background and objects while the refinement network is used to fine-tune a rough input mask.



**Fig. 3** Advantages of using the smart dilation mask, i.e. a smart border layer in the output map of our Multi-OSVOS network. (a) The border is obtained by simply dilating the output map of the network: some parts of the objects are not covered. (b) The border layer is learned by the network: the transition region is covered.

motion blur. A typical example can be seen in figure 3 where some parts of the man's hands and legs cannot be captured by simply dilating the output mask because motion blur leads to partially transparent zones which are not recognized by the network as part of the man's body. With the smart dilation mask, the missing parts are properly captured, and there are no leftover pixels.

**Refinement network** The multi-OSVOS network can separate objects and background precisely, but it relies exclusively on how they appear in the annotated frame without consideration of their position, shape or motion cues across frames. Therefore, when objects have similar appearance, multi-OSVOS fails to separate between individual object instances. In order to take such cues into account we propagate and compare the prediction of multi-OSVOS across frames using video tracking techniques (Section 3.2.2) and then we double-check and improve the result after each tracking step using the refinement network described below.

The refinement network has the same architecture as the multi-OSVOS network, except that (a) it takes an additional input, namely mask predictions for the  $K$  foreground objects from another method, and (b) it does not produce as an output the  $(K + 1)$ -th smart dilation mask that does not require any further

improvement for our purposes.

Training is performed in exactly the same way as for multi-OSVOS, except that the training set has to be augmented with inaccurate input mask predictions. These should not be exactly the same as the output masks, otherwise the network would learn to perform a trivial operation ignoring the RGB information. Such inaccurate input mask predictions are created by applying relevant random degradations to ground truth masks, e.g., small translations, affine and thin-plate spline deformations, followed by a coarsening step (morphological contour smoothing and dilation) to remove details of the object contour; finally, some random tiny square blocks are added to simulate common errors in the output of multi-OSVOS. The ground truth output masks in the training dataset are also dilated by a structuring element of size  $7 \times 7$  pixels in order to have a safety margin which ensures that the mask does not miss any part of the object.

### 3.2.2 Multiple object tracking

As a complement to CNN-based segmentation we use more classical video tracking techniques in order to take global motion and position information into account. The simplest ingredient of our object tracking subsystem is a motion-based *mask propagation* technique that uses a patch-based similarity measure to propagate a known mask to the consecutive frames.



It corresponds to block (b) in Figure 1 and it will be described in more detail below. This simple scheme alone can provide results similar to other object tracking methods such as SeamSeg [68] or ObjectFlow [75]. In particular it is able to distinguish between different instances of similar objects, based on motion and position. However it loses track of the objects when they cross each other, and it accumulates the errors. To prevent this from happening we complement the mask propagation module with five coherence reinforcement steps:

**Semantic segmentation:** The refinement network (Section 3.2.1) is applied to the output of each mask propagation step in order to avoid errors accumulating from one frame to the next.

**Keyframe extraction:** Mask propagation is effective only when it propagates from frames where object masks are accurate (especially when objects do not cross each other). Frames where this is detected to be true are labeled as *keyframes*, and mask propagation is performed only between pairs of successive keyframes.

**Mask linking:** When the mask propagation step is not sure about which decision to make, it will provide not one, but several mask candidates for each object. A graph-based technique allows to link together all these mask candidates. This way the decision on which mask candidate is the best for a given object on a given frame is taken based on global motion and appearance information.

**Post-processing:** After mask linking a series of post-processing steps are performed that use the original Multi-OSVOS result to expand labelling to unlabelled regions.

**Interactive correction:** In some situations where errors appear, the user can manually correct them on one frame and this correction is propagated to the remaining frames by the propagation module.

The following paragraphs describe in detail the inner workings of the four main modules of our multiple object tracking subsystem: (a) Keyframe extraction, (b) Mask propagation, (c) Mask linking and (d) Post-processing.

**Keyframe extraction.** A frame  $t$  is a keyframe for an object  $i \in \{1, \dots, K\}$  if the mask of this particular object is known or can be computed with high accuracy.

All frames where the object masks were manually provided by the user are considered keyframes. This is usually the first frame or very few representative frames.

The remaining frames are considered keyframes for a particular object when the object is clearly isolated from other objects and the mask for this object can be computed easily. To quantify this criterion, we rely on the multi-OSVOS network which returns  $K + 1$  masks  $O_i$  for each frame  $t$  and  $i \in \{1, \dots, K + 1\}$ . This allows to compute the global foreground mask  $F = \bigcup_{i=1}^{K+1} O_i$ . To verify if this frame is a keyframe for object  $i \in \{1, \dots, K\}$  we proceed as follows:

1. Compute the connected components of  $O_i$ . Let  $O'_i$  represent the largest connected component.
2. Compute the set of connected components of the global foreground mask  $F$  and call it  $\mathcal{F}$ .
3. For each connected component  $O' \in \mathcal{F}$  compute the overlap ratio with the current object  $r_i(O') = \frac{|O'_i \cap O'|}{|O'|}$ . If  $r_i(O') > 80\%$  and both  $O'_i$  and  $O'$  are isolated from the remaining objects<sup>3</sup> then this is a keyframe for object  $i$ .

**Mask propagation** Masks are propagated forwards and backwards between keyframes to ensure temporal coherence. More specifically, the forward propagation proceeds as follows: Given the mask  $M_t$  at frame  $t$ , the propagated mask  $M_{t+1}$  is constructed with the help of a patch-based nearest neighbor shift map  $\phi_t$  from frame  $t + 1$  to frame  $t$ , defined as

$$\phi_t(p) := \operatorname{argmin}_{\delta} \underbrace{\sum_{q \in N_p} \|u_{t+1}(q) - u_t(q + \delta)\|^2}_{d^2(D_{t+1}(p), D_t(p+\delta))}$$

*i.e.* it is the shift  $\delta$  that minimizes the squared Euclidean distance between the patch centered at pixel  $p$  in frame  $t + 1$  and the patch around  $p + \delta$  at frame  $t$ . In this expression,  $N_p$  denotes a square neighborhood of given size centered at  $p$ , and  $D_t(p)$  is the associated patch in frame  $t$ , *i.e.*  $D_t(p) = u_t(N_p)$  with  $u_t$  the RGB image corresponding to frame  $t$ . The  $\ell^2$ -metric between patches is denoted as  $d$ . To improve robustness and speed, this shift map is often computed using an approximate nearest neighbor searching algorithm such as Coherency Sensitive Hashing (CSH) [38], or FeatureMatch [67]. To capture the connectivity of patches across frames in the video, two additional terms are used in [68] for space and time consistency:

<sup>3</sup>*i.e.* if  $O'_i \cap O'_j = O' \cap O'_j = \emptyset$  for all  $j \in \{1, \dots, K\}$  such that  $j \neq i$

the first term penalizes the absolute shift and the latter penalizes neighbourhood incoherence to ensure adjacent patches flow coherently. Moreover, to reduce the patch space dimension and to speed up the search, all patches are represented with lower dimension features, e.g. the main components in the Walsh-Hadamard space, see [68] for more details. We use this model to calculate our shift map.

Once the shift map has been computed we propagate the mask as follows: Let  $u_t(p)$  be the RGB value of the pixel  $p$  in frame  $t$ , then the similarity between a patch  $D_{t+1}(p)$  in frame  $t+1$  and its nearest neighbour  $D_t(p + \phi(p))$  in frame  $t$  is measured as

$$s_p = \exp(-d^2(D_{t+1}(p), D_t(p + \phi_t(p)))).$$

Using this similarity measure the mask  $M_{t+1}$  is propagated from  $M_t$  using the following rule:

$$\tilde{M}_{t+1}(p) = \begin{cases} 1 & \text{if } \sum_{q \in N_p} s_q M_t(q) > \frac{1}{2} \sum_{q \in N_p} s_q \\ 0 & \text{otherwise} \end{cases}$$

The final propagated mask  $M_{t+1}$  is obtained by a series of morphological operations including opening and hole filling on  $\tilde{M}_{t+1}$  followed by the refinement network to correct some errors. Then  $M_{t+1}$  is iteratively propagated to the next frame  $t+2$  using the same procedure until we reach the next keyframe.

Although this mask propagation approach is useful, several artifacts may occur when objects cross each other: the propagation algorithm may lose track of an occluded object or it could mistake one object for the other.

To avoid such errors, mask propagation is performed in both forwards and backwards directions between keyframes. This gives for each object two candidate masks at each frame  $t$ :  $M_t^1 = M_t^{FW}$ , *i.e.* the one that has been forward-propagated from a previous keyframe  $t' < t$  and  $M_t^2 = M_t^{BW}$ , *i.e.* the one that has been backward-propagated from an upcoming keyframe  $t' > t$ . In order to circumvent both lost and mistaken objects we consider for each object two additional candidate masks:

$$M_t^3 = M_t^{FW} \cap M_t^{BW} \quad \text{and} \quad M_t^4 = M_t^{FW} \cup M_t^{BW}.$$

The decision between these four mask candidates for each frame and each object is deferred to the next step, which makes that decision based on a global optimization.

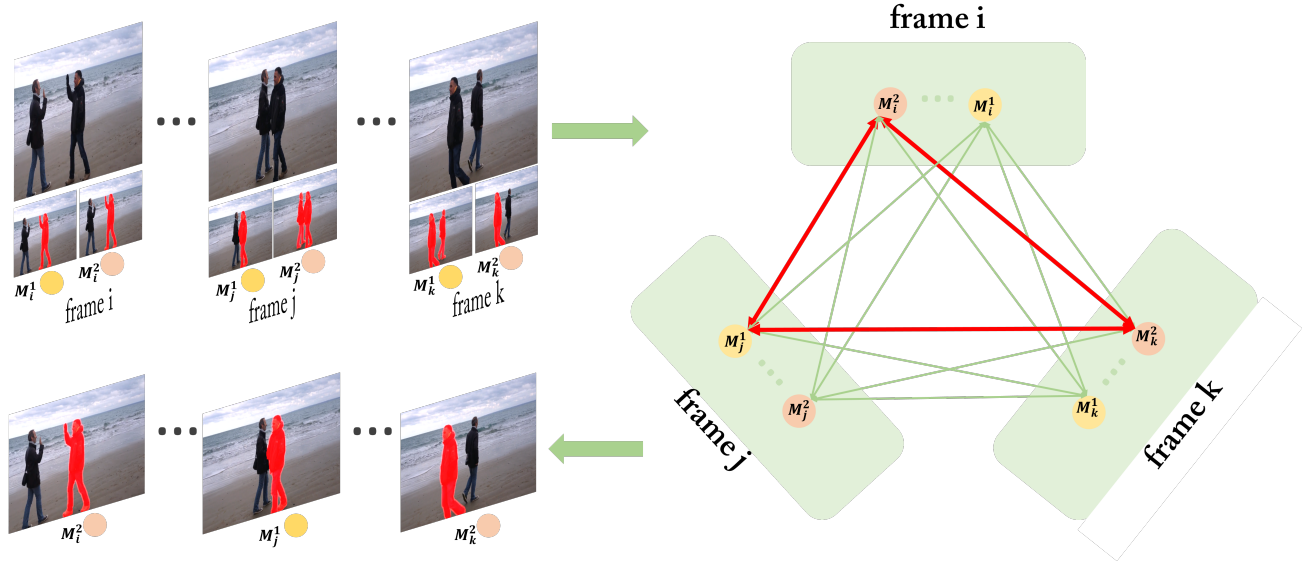
**Mask linking** After the backward and forward propagation, each object has 4 mask proposals (except for keyframes where it has a single mask proposal). In order to decide which mask to pick for each object at each frame, we use a graph-based data association

technique (GMMCP) [21] that is specially well-suited for video tracking problems. This technique does not only allow to select among the 4 candidates for a given object on a given frame. It is also capable of correcting erroneous object-mask assignments on a given frame, based on global similarity computations between mask proposals along the whole sequence. The underlying generalized maximum multi-cliques problem is clearly NP-hard, but the problem itself is of sufficiently small size to be handled effectively by a fast Binary-Integer Program as in [21].

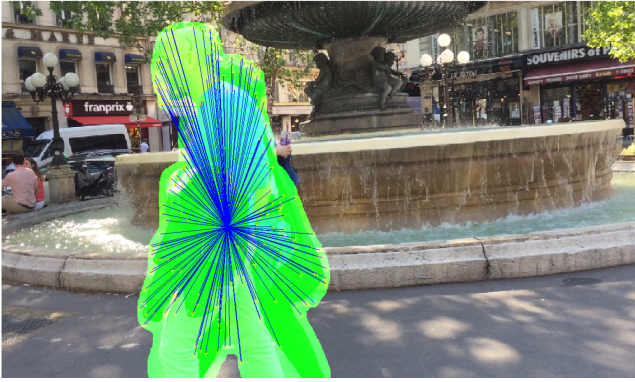
Formally, we define a complete undirected graph  $G = (V, E)$  where  $V$  is a set of vertices, each vertex corresponding to a mask proposal. Vertices in the same frame are grouped together to form a cluster.  $E$  is the set of edges connecting any two different vertices. Each edge  $e \in E$  is weighted by a score measuring the similarity between the two masks it connects. This score will be detailed in the next paragraph. All vertices in different clusters are connected together. The objective is to pick a set of  $K$  cliques<sup>4</sup> that maximize the total similarity score, with the restriction that each clique contains exactly one vertex from each cluster. Each selected clique represents the most coherent tracking of an object across all frames.

**Region similarity for mask linking** In our graph-based technique, a score needs to be specified to measure the similarity between the two masks, and the associated image data. This similarity must be robust to illumination changes, shape deformation or occlusion. Many previous approaches in multiple object tracking [21, 69] have focused on global information of the appearance model, typically the global histogram, or motion information (given by the optical flow or a simple constant velocity assumption). However, when dealing with large displacement and with an unstable camera, the constant velocity assumption is invalid and optical flow estimation is hard to apply. Furthermore, using only global information is not sufficient since our object regions already resemble in global appearance. To overcome this challenge, we define our similarity score as the combination between global and local features. More precisely, each region  $R$  is described by the corresponding mask  $M$  its global HSV histograms  $H$ , a set  $P$  of SURF keypoints [3] in it and a set  $E$  of vectors which connect each keypoint with the centroid of the mask. Each region is determined by **four elements**: \_\_\_\_\_

<sup>4</sup>A clique is a subgraph in which every two distinct vertices are connected.



**Fig. 4** Mask proposals are linked across frames to form a graph. The goal is then to select a clique from this graph minimizing the overall cost. As a result, a best candidate is picked for each frame to ensure that the same physical object is tracked.



**Fig. 5** Region description, each region is described by a global histogram, a set of SURF keypoints (yellow points), and a set of vectors which connects each keypoint and the centroid of the region.

$$R := (M, H, P, E)$$

$P := \{p_1, p_2, \dots, p_N \mid p_i \in M\}$  where  $p_i$  is the  $i$ -th keypoint

$E := \{\vec{e}_1, \vec{e}_2, \dots, \vec{e}_N \mid \vec{e}_i = p_i - C\}$  where  $C$  is the barycenter of  $M$ .

Then the similarity between two regions is defined as:

$$S(R_1, R_2) = S_H(R_1, R_2) + \alpha S_P(R_1, R_2)$$

In this expression,  $S_H(R_1, R_2) = \exp(-d_c(H_1, H_2))$  where  $d_c$  is the cosine distance between two HSV histograms which encode global color information,  $S_P$  is the local similarity computed based on keypoint matching, and  $\alpha$  is the balance coefficient to specify the contribution of each component.  $S_P$  is computed

by

$$S_P(R_1, R_2) = \sum_{p_i \in P_1} \sum_{p_j \in P_2} \gamma_{ij} \cdot w_{ij}$$

where  $\gamma_{ij}$  is the indicator function which is set to 1 if two keypoints  $p_i$  and  $p_j$  match, and zero otherwise. This function is weighted by  $w_{ij}$  based on the position of the matching keypoints with respect to the centroid of the region:

$$w_{ij} = \exp\left(\frac{-d(\vec{e}_i, \vec{e}_j)}{2\sigma}\right)$$

where  $d_c$  is the cosine distance between two vectors and  $\sigma$  is a constant.

**Post-processing** At this time, we already have  $K$  masks for  $K$  objects for all frames in video. Now we perform a post-processing step to make sure our final mask covers all the details of the objects. This is very important in video object removal since any missing detail can cause perceptually annoying artifacts in the object removal result. This post-processing includes two main steps:

The first step is to give a label for each region in the global foreground mask  $F_t = \bigcup_{i=1}^K O_t^i$  (the union of all object masks produced by multi OSVOS for frame  $t$ ) which does not have any label yet. To this end, we proceed as follows: First, we compute the connected components  $C$  of all masks  $O_t^i$  and try to assign a label to all pixels in each connected component. To this end we consider the masks  $M_t^j$  that were obtained for the same frame  $t$  (and possibly another object class  $j$  by the mask linking method). A connected

component is considered as isolated if  $C \cap M_t^j$  is empty for all  $j$ . For non-isolated components a label will be assigned by a voting scheme based on the ratio  $r_j(C) = \frac{|C \cap M_t^j|}{|C|}$ , *i.e.* the assigned label for region  $C$  will be  $\hat{j} = \operatorname{argmax}_j r_j(C)$ , the one with the highest ratio. If  $r_j(C) > 80\%$  then region  $C$  is also assigned label  $j$  regardless of the voting result, which leads to possibly multiple labels per pixel.

In the second step, we do a series of morphological operations, namely opening and hole filling. Finally we dilate each object mask again with size  $9 \times 9$ , this time allowing overlap between objects.

### 3.3 Object removal

Following the method from the previous section, all selected objects have been segmented along the complete video sequence. From the corresponding masks, the user can then decide the objects to be removed. This last step is performed thanks to video inpainting techniques that we now detail. First, we present a simple inpainting method that is adapted to the case where the background is static (or can be stabilized) and revealed at some point in the sequence. This first method is fast and relies on the reconstruction of a motion field. Then, we present a more involved method for the case where the background is moving, with possibly some complex motion as in the case of dynamic textures.

#### 3.3.1 Static background

We assume for this first inpainting method that the background is visible at least in some frames (for instance because the object to be removed is moving over a large enough distance). We also assume that the background is rigid and that its motion is only due to the camera motion. In this case, the best option to perform inpainting is to copy the visible parts of the background into the missing regions, from either past or future frames. For this, the idea is to rely on a simple optical-flow pixel propagation technique. Motion information is used to track the missing pixels and establish a trajectory from the missing region toward the source region.

**Overview of the method** Our optical flow-based pixel propagation approach is composed of three main steps, as illustrated in Figure 6. After stabilizing the video to compensate the camera movements, we use FlowNet 2.0 to estimate forward and backward optical flow fields. These optical flow fields are then inpainted using a classical image inpainting method to fill in

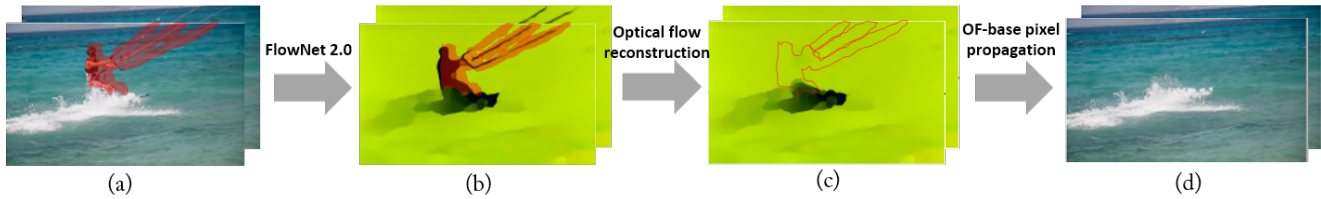
the missing information. Next, these inpainted motion fields are concatenated to create a correspondence map between pixels in the inpainting region and known pixels. Lastly, missing pixels are reconstructed by a copy-paste scheme followed by a Poisson blending to reduce artifacts.

**Motion field reconstruction** A possible approach to optical flow inpainting is smooth interpolation, for instance, in the framework of a variational approach, by ignoring the data term and using only the smoothness term in the missing regions, as proposed in [5, 88]. However, this approach leads to over-smoothed and unreliable optical flow. Therefore, we choose to reconstruct the optical flow using more sophisticated image inpainting techniques. More specifically we first compute, outside the missing region, forward/backward optical flow fields between two consecutive frames using the FlowNet approach from [34]. We then rely on the image inpainting method from [55] to interpolate these motion fields.

**Optical flow-based pixel reconstruction** Once the motion field inside the missing region is filled, it is used to propagate pixel values from the source toward the missing regions. For this to be done, we map each pixel in the missing region to a pixel in the source region. This map is obtained by accumulating the optical flow field from frame to frame (with bilinear interpolation). We do both forward and backward optical flow, which leads us to two correspondance maps: forward map and backward map. From either map, we can reconstruct missing pixels with a simple copy-paste method, using the known values outside the missing region.

We perform two passes: first a forward pass using the forward map to reconstruct the occlusion, then a backward pass using the backward map. After these two passes, the remaining missing information corresponds to parts that have never been revealed in the video. To reconstruct this information, we first use the image inpainting method from [55] to complete one keyframe, which is chosen to be the middle frame of the video, and then propagate information from this frame to other frames in the video using forward and backward maps.

**Poisson blending** Videos in real-life often contain illumination changes, especially when they are recorded outdoor. This is problematic for our approach that simply copy-paste pixel values. When the illumination of the sources is different from the illumination of the



**Fig. 6** The global pipeline of the optical flow-based propagation approach for reconstructing a static background: From input video (a), forward/backward optical flow fields are estimated by FlowNet 2.0 (b), then are inpainted by an image inpainting algorithm (c). From these optical flow fields, pixels from the source region are propagated into the missing region (d).

restored frame, visible artifacts across the border of the occlusion may appear. A common way to resolve this is by applying a blending technique, e.g. Poisson blending [65], which fuses a source image and a target image in the gradient domain. However, doing Poisson blending frame-by-frame may affect the temporal consistency. To maintain it, we adopt the recent method of Bokov *et al.* [5] which takes into account the information of the previous frame. In this method, a regularizer which penalizes discrepancies between the reconstructed colors and their corresponding colors in the optical-flow-aligned previous frame is introduced. More specifically, given the colors of the current and previous inpainted frames  $I_t(p)$ ,  $I_{t-1}(p)$ , respectively, the refined Poisson-blended image  $I(p)$  can be obtained by minimizing the discretized energy functional [5]:

$$\begin{aligned}
 B(I) = & \sum_{p \in \Omega_t} \|\nabla I(p) - G_t(p)\|^2 \\
 & + \sum_{p \in \partial\Omega_t} w_p^{PB} \|I(p) - I_t(p)\|^2 \\
 & + \sum_{p \in \Omega_t} (1 - w_p^{PB}) \|I(p) - I_{t-1}(p + O_t(p))\|^2
 \end{aligned}$$

Here,  $\partial\Omega_t$  denotes the outer-boundary pixels of the missing region  $\Omega_t$ ,  $G_t(p)$  is the target gradient field and  $O_t(p)$  is the optical flow at position  $p$  between frames  $t-1$  and  $t$ . The terms  $w_p^{PB}$  are defined as

$$w_p^{PB} = (1 + \sigma^{PB} \|\nabla I^{PB}(p) - G_t(p)\|^2)^{-1},$$

where  $I^{PB}$  is the usual Poisson blended image, and are used to weight the reconstruction results from the previous frame  $I_{t-1}$  in the boundary conditions. In this definition,  $\sigma^{PB}$  is a constant controlling the strength of the temporal-consistency enforcement. These weights allow to better deal with global illumination changes while enforcing temporal stability. This Poisson blending technique is applied at every pixel propagation step to support the copy-paste framework.

### 3.3.2 Dynamic background

The simple optical flow-based pixel propagation method that we proposed in section 3.3.1 can

produce plausible results if the video contains only static background and simple camera motion. More involved methods are needed to deal with large pixel displacement and complex camera movements. They are typically based on joint estimation of optical flow and color information inside the occlusion, see for instance [32, 81]. However, when the background is dynamic or contains moving objects, these latter methods often fail to capture oscillatory patterns in the background. In that situation, global patch-based methods are preferred. They rely on the minimization of a global energy computed over space-time patches. This idea was first proposed in [79], later improved in [54], and recently improved again in Le *et al.* [39].

Let us describe briefly the method proposed in [39]. A prior stabilization process is applied to compensate the instabilities due to camera movements (see below for the improvement proposed in the current work). Then a multiscale coarse-to-fine scheme is used to compute a solution to the inpainting problem. The general structure of this scheme is the following: at each scale of a multiscale pyramid, we alternate until convergence the computation of an optimal shift map between pixels in the inpainting domain and pixels outside (using a metric between patches which involves image colors, texture features, and optical flow), and the update of image colors inside the inpainting domain (using a weighted average of the values provided by the shift map). A key to the quality of the final result is the coarse initialization of this scheme; it is obtained by progressively filling in the inpainting domain (at the coarsest scale) using patch matching and (mapped) neighbors averaging together with a priority term based on optical flow. The heavy use of optical flow at each scale helps a lot to enforce the temporal consistency even in difficult cases such as dynamic background or complex motion. In particular, the method can reconstruct moving objects even when they interact with each other. The whole method is computationally heavy but the speed is significantly boosted when all

steps are parallelized.

We have recently brought several improvements to this method of [39]:

**Video stabilization** : In general, patch-based video inpainting techniques require a good video stabilization as a pre-processing step to compensate patch deformations due to camera motions [56, 70]. This video stabilization is usually done by calculating a homography between two consecutive frames using keypoints matching followed by a RANSAC algorithm to remove outliers [15]. However, large moving objects appearing in the video may reduce the performances of such an approach because too many keypoints may be selected on these objects and prevent the homography from being estimated accurately from the background. This problem can be solved by simply neglecting all segmented objects for computing the homography. This is easy to do: since we already have the masks of the selected objects, we just have to remove all keypoints which are covered by masks. This is an advantage of our approach where both segmentation and inpainting are addressed.

**Background/Foreground inpainting** : In addition to stabilization improvement, multiple segmentation masks are also helpful for inpainting separately the background and the foreground. More precisely, we first inpaint the background neglecting all pixels contained in segmented objects. After that, we inpaint in priority the segmented objects that we want to keep and which are partially occluded. This increases the quality of the reconstruction, both for the background and for the objects. Furthermore, it reduces the risk of blending segmented objects which are partially occluded because segmented objects have separate labels. In particular, it is extremely helpful when several objects overlap.

Let us finally mention another advantage of our joint tracking/inpainting method: objects are better segmented and thus easier to inpaint for it is a well-known fact that the inpainting of a missing domain may be of lower quality if the boundary values are not suitable. In our case, time continuity of segmented objects and the fact of using different labels for different objects have a huge impact on the quality of the inpainting.

## 4 Results

We first evaluate our results for the segmentation step of the proposed method, for which we provide quantitative and visual results and comparisons with

state-of-the-art methods. We then provide several visual results for the complete object removal process, again comparing with the most efficient methods. These visual comparisons are given as isolated frames in the paper and it is of course more informative to go for the complete videos in the supplementary material, see <https://object-removal.telecom-paristech.fr/>. We consider various datasets: we use sequences from the DAVIS-2016 [63] challenge, from the MOViC [14], and from the ObMIC [86] datasets; we also consider classical sequences from the papers [26] and [54]. Eventually, we provide several new challenging sequences containing strong appearance changes, motion blur, objects with similar appearance and possibly crossing, as well as complex dynamic textures.

Concerning the number of annotated frames: Unless otherwise stated only the first frame is annotated by the user in all experiments. In some examples (*e.g.* CAMEL) not all objects are visible in the first frame and we use another frame for annotations. In a few examples we annotate more than one frame (*e.g.* first and last frame in TEDDY BEAR-FIRE AND JUMPING GIRL-FIRE) in order to illustrate the flexibility of the system for correcting errors.

### 4.1 Implementation details

For the segmentation part, we use the Deeplab v2 [10] architecture for the multi-OSVOS and refining networks. We initialize the network using the pre-trained model provided by [10] and then adapt it to video using the training set of DAVIS-2016 [64] and *train-val* set in DAVIS-2017 [66] datasets (sequences from which we exclude the validation set of DAVIS-2016). For the data augmentation procedure, we generate 100 pairs of images and ground truth from the first frame annotation, following the same protocol as in [37]. For the patch-based mask propagation and mask linking, we evolved from the implementation of [68] and [21], respectively.

For the video inpainting step, we use the default parameters from our previous work [39]. In particular, the patch size is set to 5, and the number of levels in the multi-scale pyramid is 4.

For a typical sequence with resolution ( $854 \times 480$ ) and 100 frames, the full computational time is of the order of 45 minutes for segmentation plus 40 minutes for inpainting on a core I7 CPU machine with 32 Gb of RAM and a GTX 1080 GPU. While this is a limitation of the approach, the complete object removal is about one order of magnitude faster than the single

completion step from state-of-the-art methods [54] or [32]. While interactive editing is out of reach for now, these computational times allow the offline post-processing of sequences.

## 4.2 Object segmentation

For the proposed object removal system, and as explained in detail above, the most crucial point is that the segmentation masks shall completely cover the considered objects, including motion and transition blur. Otherwise, unacceptable artifacts remain after the full object removal procedure (see Figure 13 for an example). In terms of performance evaluation, this means that we favor *recall* over *precision*, as defined below. This also means that the ground truth provided with classical datasets may not be fully adequate to evaluate segmentation in the context of object removal, because they do not include transition zones induced by, e.g., motion blur. For this reason, recent video inpainting methods that make use of these databases to avoid the tedious manual selection of objects, usually start from a *dilation* of the ground truth. In our case, a dilation is learned by our architecture (smart dilation) at the segmentation step, as explained above. For these reasons, we compare our method with state-of-the-art object segmentation methods, after various dilations and on the dilated versions of the ground truth. We also provide visual results in our supplementary website: <https://object-removal.telecom-paristech.fr/>.

**Evaluation Metrics** We briefly recall here the evaluation metrics that we use in this work: some of them are the same as in the DAVIS-2016 challenge [63] and we also add other metrics that are specialized for our task. The goal is to compare the computed segmentation mask (SM) to the ground truth mask (GT). The *recall* metric is defined as the ratio between the area of the intersection between SM and GT, and the area of GT. The *precision* is the ratio between the area of the intersection and the area of the SM. Eventually, the *IOU* (intersection over union), or Jaccard index, is defined as the ratio between intersection and union.

**Single object segmentation** We use the DAVIS-2016 [63] validation set and compare our approach to recent semi-supervised state-of-the-art techniques (SeamSeg [68], ObjectFlow [75], MSK [64], OSVOS [8] and onAVOS [77]) using the pre-computed segmentation masks provided by the authors. As explained above, we consider a dilated version of the

	Metric		
	Recall (%)	Precision (%)	IoU (%)
SeamSeg	59.31	73.08	50.20
ObjectFlow	70.63	90.97	67.78
MSK	82.83	95.00	79.94
OSVOS	86.78	92.38	80.58
onAVOS	87.64	<b>96.67</b>	<b>85.17</b>
Ours	<b>89.63</b>	94.31	84.70

**Tab. 1** Quantitative evaluation of our object segmentation method compared to other state-of-the-art methods, on the single object DAVIS-2016[63] validation set. As explained in the text, the main objective when performing object removal is to achieve high Recall scores.

ground truth (we use a dilation by a  $15 \times 15$  structuring element, as in [32, 39]). Therefore, we apply a dilation of the same size to the masks from all the concurrent methods. In our case, this dilation has both been learned (size  $7 \times 7$ ) and applied as a post-processing step (size  $9 \times 9$ ). Since the composition of two dilations with such sizes yields a dilation with size  $15 \times 15$ , the comparison is fair.

Table 1 shows the comparisons using the three above-mentioned metrics. Our method has the best recall score overall, therefore achieving its objective. The precision score remains very competitive. Besides, our method outperforms OSVOS [8] and MSK [64], those having a similar neural network backbone architecture (VGG16), on all metrics. The precision and *IOU* scores compare favorably with onAVOS [77] which uses a deeper and more advanced network. Table 2 provides a comparison between OSVOS [8] and our approach on two sequences from [27]. These sequences have been manually segmented by the authors of [27] for video inpainting purposes. On such extremely conservative segmentation masks (in the sense that they over-detect the object), the advantage of our method is particularly strong.

As a further experiment, we investigate the ability of dilations with various sizes to improve the recall without degrading the precision too much. For this, we plot precision-recall curves as a function of the structuring element size (ranging from 1 to 30). To include our method on this graph, we start from our original method (highlighted with a green square) and apply to it either erosions with a radius ranging from 1 to 15, or dilation with a radius ranging from 1 to 15. Again this makes sense since our method has learned a dilation whose equivalent radius is 15. Results are displayed in Figure 8. As can be seen from this

	Metric		
	Recall (%)	Precision (%)	IoU (%)
<b>Granados-S1</b>			
OSVOS	62.04	59.17	52.15
Ours	80.12	86.31	67.53
<b>Granados-S3</b>			
OSVOS	74.42	87.00	63.02
Ours	80.12	86.31	67.53

**Tab. 2** Quantitative evaluation of our object segmentation method compared to the OSVOS segmentation method [8], on two sequences manually segmented for inpainting purposes [27]

	Metric		
	Recall (%)	Precision (%)	IoU (%)
<b>MOVICs</b>			
SeamSeg	78.63	74.06	65.96
ObjectFlow	59.50	77.01	52.33
OSVOS	85.48	83.87	76.63
Ours	<b>89.28</b>	<b>87.09</b>	<b>81.58</b>
<b>ObMIC</b>			
SeamSeg	91.00	80.30	75.33
ObjectFlow	53.14	83.00	43.64
OSVOS	85.89	84.08	74.55
Ours	<b>94.42</b>	<b>88.48</b>	<b>83.81</b>

**Tab. 3** Quantitative evaluation of our object segmentation method compared to other state-of-the-art methods, on two multiple objects datasets (MOVICs [14] and ObMIC [86])

figure, our method is the best in terms of recall, and the recall is increasing significantly with respect to the dilation size. With the sophisticated onAVOS method, on the other hand, the recall increases slowly, and the precision drops drastically as the dilation size increases. Basically, these experiments show that the performances achieved by our system for the full coverage of a single object (that is, with as few missed pixels as possible) cannot be obtained from state-of-the-art object segmentation methods by using simple dilation techniques.

**Multiple objects segmentation** Next, we perform the same experiments for datasets containing videos with multiple objects. Since the *test* ground truth was not yet available (at the time of this writing) for the DAVIS-2017 dataset and since our network was trained on the *train-val* set of this dataset, we consider two other datasets: MOVICs [14] and ObMIC [86]. The datasets include multiple objects, but only have one label per sequence. To evaluate the multiple

object situations, we only kept sequences containing more than one object, and then manually re-annotated the ground truth giving different labels for different instances. Observe that these datasets contain several major difficulties such as large camera displacement, motion blur, similar appearances, and crossing objects. Results are summarized in Table 3. From this table, roughly the same conclusions as in the single object can be drawn, namely the superiority of our method in term of recall score, without sacrificing much the precision score.

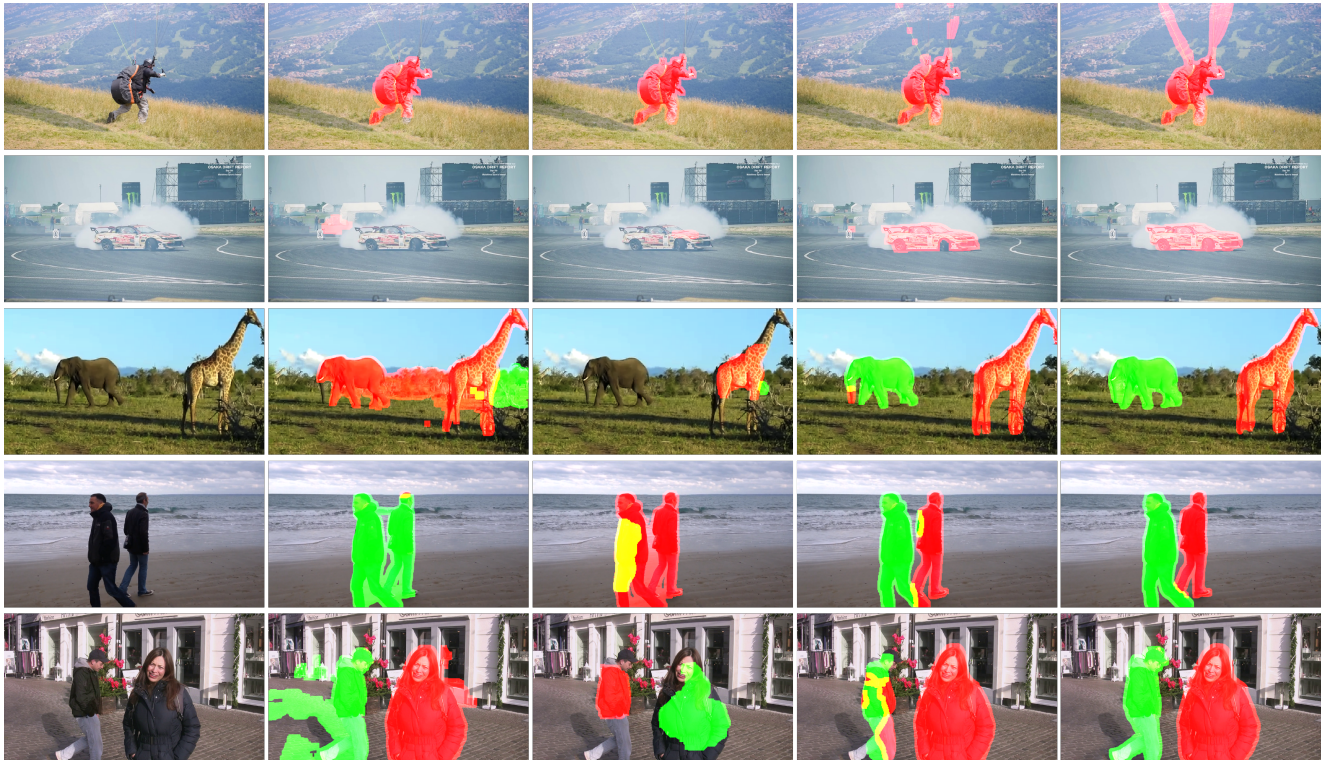
Some qualitative results of our video segmentation technique are shown in Figure 7. In the first two rows, we show some frames corresponding to the single object case, on the DAVIS-2016 dataset [63]. The last three rows show multiple objects segmentation results on MOVICs [14], ObMIC [86] and Granados's sequences [26] respectively. We observe on these examples that our approach yields full object coverage, even with complex motion and motion blur. This is particularly noticeable on the sequences KITE-SURF and PARAGLIDING-LAUNCH. In the multiple objects cases, the examples illustrate the capacity of our method to deal with complex occlusions. This cannot be achieved with mask tracking methods such as objectFlow [75] or SeamSeg [68]. The OSVOS method [8] yields some confusion between objects, probably because the temporal continuity is not taken into account by this approach.

### 4.3 Object removal

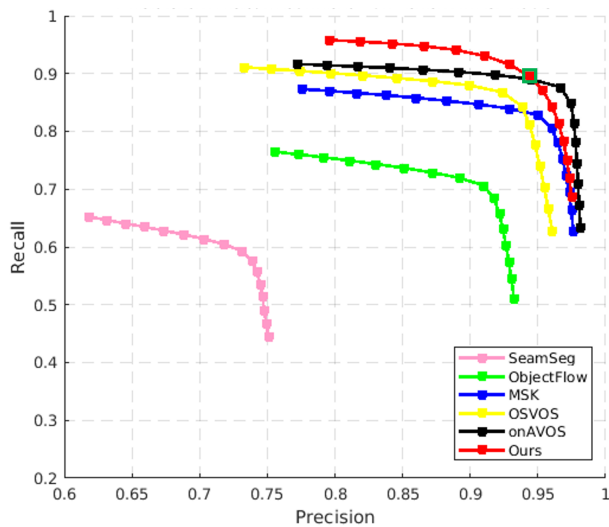
Next, we evaluate the complete object removal pipeline. We consider both the inpainting versions that we have introduced. We use the simple, optical-flow based method introduced in Section 3.3.1 for sequence having static background. We refer to this fast method as the *static version*. We use the more complex method derived from [39] and detailed in Section 3.3.2 for more involved sequences, exhibiting challenging situations such as dynamic background, camera instability, complex motions, and crossing objects. We refer to this second slower version as the *dynamic version*.

In Figure 9, we display examples of both single and multiple objects removal, through several representative frames. The video results can be fully viewed in the supplementary website. The first sequence BLACKSWAN (DAVIS-2016) shows that our method (dynamic version) can plausibly reproduce dynamic textures. In the second sequence COWS (DAVIS 2016), the method yields good results, with





**Fig. 7** Visual comparison of different segmentation approaches. From left to right, Original, SeamSeg [68], ObjectFlow [75], OSVOS [8], ours.



**Fig. 8** Precision-recall curves for different methods with different dilation sizes.

a stable background and continuity of the geometrical structures, despite a large occlusion implying that some regions are covered through all the sequence. We then turn to the case of multiple objects removal. In the sequence CAMEL (DAVIS-2017), we show the removal of one static object, a challenging case since the background information is missing at places. On this example, the direct use of the inpainting method from [39] results in some undesired artifacts when the second camel enters the occlusion. By using multiple object segmentation masks to separate background and foreground, we can create a much more stable background. The last two examples are from an original video. This sequence again highlights that our method can deal with dynamic textures and hand-held cameras.

**Comparison with state-of-the-art inpainting methods** In these experiments, we compare our results with the state-of-the-art video inpainting methods [32] and [54].

First, we provide a visual comparison between our optical flow-based pixel propagation (that is, the static approach) with the method of Huang *et al.* [32] using a video with a static background. Figure 10 shows some representative frames of the sequence HORSE-JUMP-



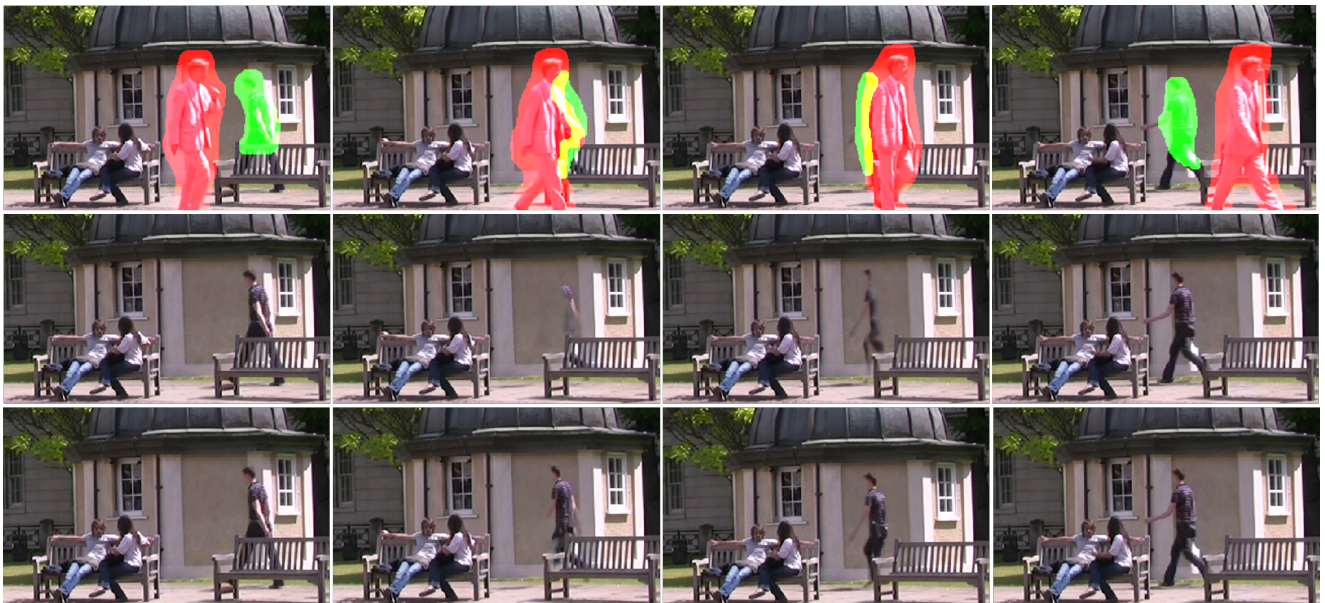
Fig. 9 Visual illustrations of our objects removal system.



**Fig. 10** Qualitative comparison with Huang's method [32]. From top to bottom: our segmentation mask, result from [32] performed on manually segmented mask, our inpainting results performed on our mask.



**Fig. 11** Qualitative comparison with Huang's method [32] on video with dynamic background. From left to right: our segmentation mask, result from [32], our inpainting result performed on our mask.



**Fig. 12** Qualitative comparison with Newson et al's method [54]. Top: our segmentation masks, red and green masks denote different objects, yellow region is the overlap region between two objects. Middle: results from [54] performed on our segmentation masks. Bottom: our inpainting results performed on the same masks.

HIGH. In this sequence, we get a comparable result using our simple optical flow-based pixel propagation approach. Our advantage is the considerable reduction of the computational time. With a not-optimized version of the code, our method takes approximately 30 minutes to finish while [32] takes about 3 hours to complete this sequence.

Next, we qualitatively compare our method with [32] when reconstructing dynamic backgrounds. We use the code released by the author on several sequences using the default parameters. In general, Huang *et al.* [32] fail to generate convincing dynamic textures. This can be explained by the fact that their algorithm relies on dense flow fields to guide the completion, these fields being often unreliable for dynamic texture. Moreover, they fill the hole by sampling only 2D patches from the source regions and therefore the periodic repetition of the background is not captured. Our method, on the other hand, fills the missing dynamic textures in a plausible way. Figure 10 shows the representative frames of the reconstructed sequence TEDDY-BEAR, which is recorded indoor. This sequence is especially challenging because of the presence of both dynamic and static textures, as well as because of illumination changes. Our method yields a convincing reconstruction of the fire, contrarily to [32]. The complete video can be seen in the supplemental material website.

We also compare our results with the video inpainting technique from [54]. Figure 12 shows some representative frames of the sequence PARK-COMPLEX, which is taken from [27] and is modified to focus on the moment where objects occlude each other. In this example, the method of [54] cannot reconstruct the moving man on the right which is occluded by the man on the left. This is because the background behind this man changes over time (from tree to wall). Since Newson *et al.*'s method [54] treats the background and the foreground similarly, the algorithm can not reconstruct the situation "man in front of the wall" because it never sees this situation before. Our method, by making use of the optical flow and thanks to the objects segmentation map, can reconstruct the "man" and the "wall" independently, yielding a plausible reconstruction.

**Impact of the segmentation masks on the inpainting performances.** In these experiments, we highlight the advantages of using the segmentation masks of multiple objects to improve the video inpainting results.

First, we emphasize the need for masks which fully

cover the objects to be removed. Figure 13 (top) demonstrate the situation where some object details (the waving hand in this case) are not covered by the mask (here using the state-of-the-art OSVOS method) [8]. This situation leads to a very unpleasant artifact when video inpainting is performed. Thanks to the smart dilation, introduced in the previous sections, our segmentation mask fully cover the object to be removed, yielding a more plausible video after the inpainting step.

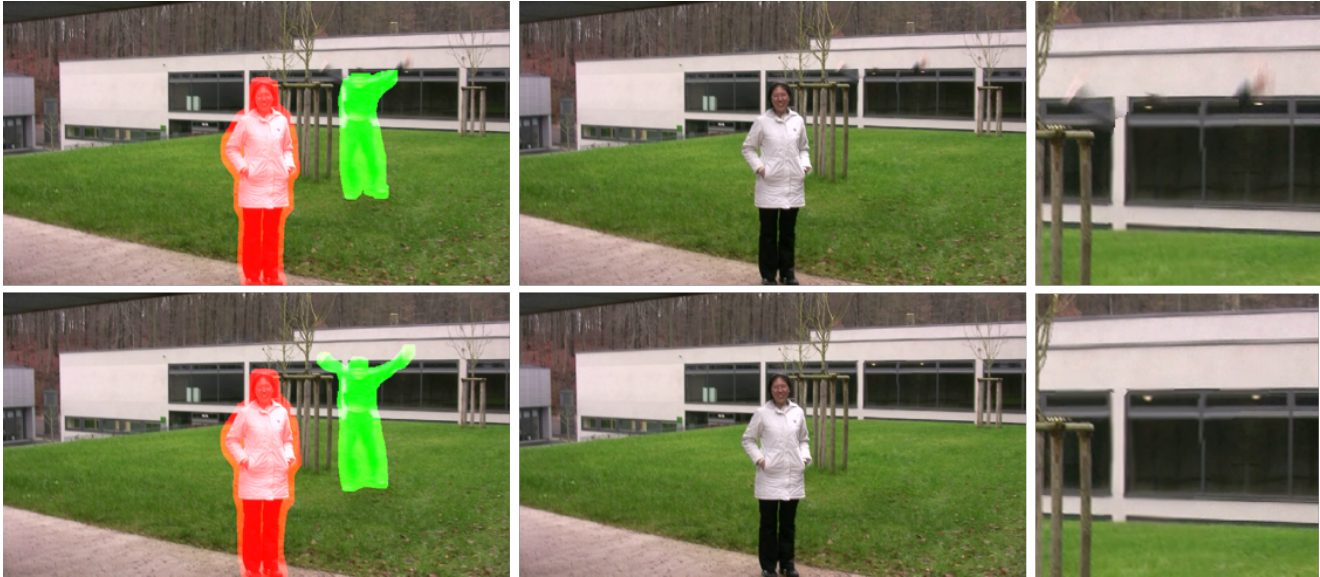
Object segmentation masks can also be helpful for the video stabilization step. Indeed, in case of large foregrounds, these can have a strong effect on the stabilization procedure, yielding a bad stabilization of the background, which in turn yields bad inpainting results. In contrast, if the stabilization is applied only to the background, the final object removal results are much better. This situation is illustrated in the supplementary material.

To further investigate the advantage of using multiple segmentation masks to separate background/foreground in the video completion algorithm, we compare our method with the direct application of the inpainting method from [39], without separating objects and background. Representative frames of both approaches are shown in Figure 14. Clearly, [39] produce artifacts when the moving objects (the two characters) overlap the occlusion, due to patches from these moving objects being propagated within the occlusion in the nearest neighbor search step. Our method, on the other hand, does not suffer from this problem because we reconstruct background and moving objects separately. This way, the background is more stable, and the moving objects are well reconstructed.

## 5 Conclusion, limitations and discussion

In this paper, we have provided a full system performing object removal in videos. The input of the system is made of a few strokes provided by the user to indicate the objects to be removed. To the best of our knowledge, this is the first system of this kind, even though the Adobe company has recently announced to be developing such a tool, under the name *Cloak*. The approach can deal with multiple, possibly crossing objects, and can reproduce complex motions and dynamic textures.

Although our method achieves good visual results on different datasets, it still suffers from a few limitations.



**Fig. 13** Results of object removal using masks computed by OSVOS (top) and ours (bottom). From left to right: Segmentation mask, the resulting object removal on one frame, zooms. We can see that when the segmentation masks do not fully cover the object (OSVOS), the resulting video contain visible artifacts (the hand of the man remains after object removal).



**Fig. 14** The advantage of using the segmentation masks to separate background and foreground. Left: without separating background/foreground, the result have many artifacts. Right: the background and foreground are well reconstructed when being reconstructed independently.

First, parts of the objects to be edited may be ignored by the segmentation masks. In such cases, as already emphasized, the inpainting step of the algorithm will amplify the remaining parts, creating strong artifacts. This is an intrinsic problem of the semi-supervised object removal task and room remains for further improvement. Further, the system is still relatively slow, and in any case far from realtime. Accelerating the system could allow for interactive scenarios where the user can gradually correct the segmentation-inpainting loop.

The segmentation of shadows is still not flawlessly performed by our system, especially when the shadows are not strongly contrasted. It is a desirable property of the system to be able to deal with such cases. This

problem can be seen in several examples provided in the supplementary material.

Concerning the inpainting module the user has to currently choose between the fast motion-based version (which works better for static backgrounds) and the slower patch-based version which is required in the presence of complex dynamic backgrounds. An integrated method that reunites the advantages of both would be preferable. Huang's method [32] makes a nice attempt in this direction, but its use of 2D patches is not sufficient to correctly inpaint complex dynamic textures, which are more plausibly inpainted by our 3D patch-based method.

Another limitation occurs in some cases where the background is not revealed, specifically when some

semantic information should be used. Such difficult cases are gradually being solved for single images by using CNN-based inpainting schemes [33]. While the training step of such methods is still out of reach for videos as of today, developing an object removal scheme fully relying on neural networks is an exciting research direction.

## Acknowledgements

We gratefully acknowledge the support of NVIDIA with the donation of the Titan Xp GPU used for this research. This work was funded by the French Research Agency (ANR) under Grant No ANR-14-CE27-001 (MIRIAM).

**Open Access** This article is distributed under the terms of the Creative Commons Attribution License which permits any use, distribution, and reproduction in any medium, provided the original author(s) and the source are credited.

## References

- [1] Y. Aksoy, T.-H. Oh, S. Paris, M. Pollefeys, and W. Matusik. Semantic soft segmentation. *ACM Trans. Graph. (Proc. SIGGRAPH)*, 37(4):72:1–72:13, 2018.
- [2] X. Bai, J. Wang, D. Simons, and G. Sapiro. Video snapcut: robust video object cutout using localized classifiers. In *ACM Transactions on Graphics (ToG)*, volume 28, page 70. ACM, 2009.
- [3] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool. Speeded-Up Robust Features (SURF). *Computer Vision and Image Understanding*, 110(3):346–359, jun 2008.
- [4] M. Bertalmio, G. Sapiro, V. Caselles, and C. Ballester. Image inpainting. In *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '00*, pages 417–424, 2000.
- [5] A. Bokov and D. Vatolin. 100+ times faster video completion by optical-flow-guided variational refinement. *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 2122–2126, 2018.
- [6] N. Bonneel, K. Sunkavalli, J. Tompkin, D. Sun, S. Paris, and H. Pfister. Interactive intrinsic video editing. *ACM Transactions on Graphics (TOG)*, 33(6):197, 2014.
- [7] S. Caelles, Y. Chen, J. Pont-Tuset, and L. Van Gool. Semantically-guided video object segmentation. *arXiv preprint arXiv:1704.01926*, 2017.
- [8] S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool. One-shot video object segmentation. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [9] S. Caelles, A. Montes, K.-K. Maninis, Y. Chen, L. Van Gool, F. Perazzi, and J. Pont-Tuset. The 2018 davis challenge on video object segmentation. *arXiv:1803.00557*, 2018.
- [10] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv preprint arXiv:1606.00915*, 2016.
- [11] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. *arXiv preprint arXiv:1802.02611*, 2018.
- [12] Y. Chen, J. Pont-Tuset, A. Montes, and L. Van Gool. Blazingly fast video object segmentation with pixel-wise metric learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1189–1198, 2018.
- [13] J. Cheng, Y.-H. Tsai, W.-C. Hung, S. Wang, and M.-H. Yang. Fast and accurate online video object segmentation via tracking parts. *arXiv preprint arXiv:1806.02323*, 2018.
- [14] W.-C. Chiu and M. Fritz. Multi-class video co-segmentation with a generative multi-video model. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 321–328. IEEE, 2013.
- [15] S. Choi, T. Kim, and W. Yu. Robust video stabilization to outlier motion using adaptive ransac. In *Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on*, pages 1897–1902. IEEE, 2009.
- [16] Y.-Y. Chuang, A. Agarwala, B. Curless, D. H. Salesin, and R. Szeliski. Video matting of complex scenes. *ACM Transactions on Graphics (ToG)*, 21(3):243–248, 2002.
- [17] A. Colombari, A. Fusiello, and V. Murino. Segmentation and tracking of multiple video objects. *Pattern Recognition*, 40(4):1307–1317, 2007.
- [18] A. Criminisi, P. Pérez, and K. Toyama. Region filling and object removal by exemplar-based image inpainting. *IEEE Transactions on image processing*, 13(9):1200–1212, 2004.
- [19] Z. Cui, O. Wang, P. Tan, and J. Wang. Time slice video synthesis by robust video alignment. *ACM Transactions on Graphics (TOG)*, 36(4):131, 2017.
- [20] J. Dai, K. He, and J. Sun. Instance-Aware Semantic Segmentation via Multi-task Network Cascades. In *(CVPR) IEEE Conference on Computer Vision and Pattern Recognition*, pages 3150–3158. IEEE, jun 2016.
- [21] A. Dehghan, S. Modiri Assari, and M. Shah. Gmmcp tracker: Globally optimal generalized maximum multi clique problem for multiple object tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4091–4099, 2015.
- [22] B. Drayer and T. Brox. Object detection, tracking, and motion segmentation for object-level video segmentation. *arXiv preprint arXiv:1608.03066*, 2016.
- [23] I. Drori, D. Cohen-Or, and H. Yeshurun. Fragment-based image completion. In *ACM Transactions on*

- graphics (TOG)*, volume 22, pages 303–312. ACM, 2003.
- [24] A. A. Efros and T. K. Leung. Texture synthesis by non-parametric sampling. In *Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*, volume 2, pages 1033–1038. IEEE, 1999.
- [25] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.
- [26] M. Granados, K. I. Kim, J. Tompkin, J. Kautz, and C. Theobalt. Background inpainting for videos with dynamic objects and a free-moving camera. In *European Conference on Computer Vision*, pages 682–695. Springer, 2012.
- [27] M. Granados, J. Tompkin, K. Kim, O. Grau, J. Kautz, and C. Theobalt. How not to be seen object removal from videos of crowded scenes. In *Computer Graphics Forum*, volume 31, pages 219–228. Wiley Online Library, 2012.
- [28] H. Grossauer. Inpainting of movies using optical flow. In *Mathematical Models for Registration and Applications to Medical Imaging*, pages 151–162. Springer, 2006.
- [29] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2980–2988. IEEE, 2017.
- [30] J. Herling and W. Broll. Pixmix: A real-time approach to high-quality diminished reality. In *Mixed and Augmented Reality (ISMAR), 2012 IEEE International Symposium on*, pages 141–150. IEEE, 2012.
- [31] Y.-T. Hu, J.-B. Huang, and A. Schwing. Maskrcnn: Instance level video object segmentation. In *Advances in Neural Information Processing Systems*, pages 324–333, 2017.
- [32] J.-B. Huang, S. B. Kang, N. Ahuja, and J. Kopf. Temporally coherent completion of dynamic video. *ACM Transactions on Graphics (TOG)*, 35(6):196, 2016.
- [33] S. Iizuka, E. Simo-Serra, and H. Ishikawa. Globally and locally consistent image completion. *ACM Transactions on Graphics (TOG)*, 36(4):107, 2017.
- [34] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *2017 IEEE conference on computer vision and pattern recognition (CVPR)*, pages 1647–1655. IEEE, 2017.
- [35] S. D. Jain and K. Grauman. Click carving: Segmenting objects in video with point clicks. *arXiv preprint arXiv:1607.01115*, 2016.
- [36] J. Jia, Y.-W. Tai, T.-P. Wu, and C.-K. Tang. Video repairing under variable illumination using cyclic motions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(5):832–839, 2006.
- [37] A. Khoreva, R. Benenson, E. Ilg, T. Brox, and B. Schiele. Lucid data dreaming for object tracking. *arXiv preprint arXiv:1703.09554*, 2017.
- [38] S. Korman and S. Avidan. Coherency sensitive hashing. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 1607–1614. IEEE, 2011.
- [39] T. Le, A. Almansa, Y. Gousseau, and S. Masnou. Motion-consistent video inpainting. In *ICIP 2017: IEEE International Conference on Image Processing*, 2017.
- [40] T. T. Le, A. Almansa, Y. Gousseau, and S. Masnou. Removing objects from videos with a few strokes. In *SIGGRAPH Asia 2018 Technical Briefs*, page 22. ACM, 2018.
- [41] M. Leake, A. Davis, A. Truong, and M. Agrawala. Computational video editing for dialogue-driven scenes. *ACM Transactions on Graphics (TOG)*, 36(130), 2017.
- [42] Y. J. Lee, J. Kim, and K. Grauman. Key-segments for video object segmentation. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 1995–2002. IEEE, 2011.
- [43] A. Levin, D. Lischinski, and Y. Weiss. A closed-form solution to natural image matting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(2):228–242, 2008.
- [44] E. Levinkov, J. Tompkin, N. Bonneel, S. Kirchhoff, B. Andres, and H. Pfister. Interactive multicut video segmentation. In *Pacific Graphics*, 2016.
- [45] F. Li, T. Kim, A. Humayun, D. Tsai, and J. M. Rehg. Video segmentation by tracking many figure-ground segments. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2192–2199, 2013.
- [46] X. Li, Y. Qi, Z. Wang, K. Chen, Z. Liu, J. Shi, P. Luo, C. C. Loy, X. Tang, A. Khoreva, et al. Video object segmentation with re-identification. In *The 2017 DAVIS Challenge on Video Object Segmentation-CVPR Workshops*, 2017.
- [47] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.
- [48] J. Luiten, P. Voigtlaender, and B. Leibe. Premvos: Proposal-generation, refinement and merging for video object segmentation. *arXiv preprint arXiv:1807.09190*, 2018.
- [49] N. Märki, F. Perazzi, O. Wang, and A. Sorkine-Hornung. Bilateral space video segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 743–751, 2016.
- [50] S. Masnou and J.-M. Morel. Level lines based disocclusion. In *Image Processing, 1998. ICIP 98. Proceedings. 1998 International Conference on*, pages 259–263. IEEE, 1998.
- [51] Y. Matsushita, E. Ofek, W. Ge, X. Tang, and H.-Y. Shum. Full-frame video stabilization with motion

- inpainting. *IEEE Transactions on pattern analysis and Machine Intelligence*, 28(7):1150–1163, 2006.
- [52] F. Meyer. Topographic distance and watershed lines. *Signal Processing*, 38(1):113–125, 1994.
- [53] N. S. Nagaraja, F. R. Schmidt, and T. Brox. Video segmentation with just a few strokes. In *ICCV*, pages 3235–3243, 2015.
- [54] A. Newson, A. Almansa, M. Fradet, Y. Gousseau, and P. Pérez. Video inpainting of complex scenes. *SIAM Journal on Imaging Sciences*, 7(4):1993–2019, 2014.
- [55] A. Newson, A. Almansa, Y. Gousseau, and P. Pérez. Non-local patch-based image inpainting. *Image Processing On Line*, 7:373–385, 2017.
- [56] J. Odobez and P. Bouthemy. Robust Multiresolution Estimation of Parametric Motion Models. *Journal of Visual Communication and Image Representation*, 6(4):348–365, dec 1995.
- [57] S. W. Oh, J.-Y. Lee, K. Sunkavalli, and S. J. Kim. Fast video object segmentation by reference-guided mask propagation. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7376–7385. IEEE, 2018.
- [58] A. Papazoglou and V. Ferrari. Fast object segmentation in unconstrained video. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1777–1784, 2013.
- [59] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2536–2544, 2016.
- [60] K. A. Patwardhan, G. Sapiro, and M. Bertalmio. Video inpainting of occluding and occluded objects. In *Image Processing, 2005. ICIP 2005. IEEE International Conference on*, volume 2, pages II–69. IEEE, 2005.
- [61] K. A. Patwardhan, G. Sapiro, and M. Bertalmio. Video inpainting under constrained camera motion. *IEEE Transactions on Image Processing*, 16(2):545–553, 2007.
- [62] F. Perazzi, A. Khoreva, R. Benenson, B. Schiele, and A. Sorkine-Hornung. Learning video object segmentation from static images. In *Computer Vision and Pattern Recognition*, 2017.
- [63] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 724–732, 2016.
- [64] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Computer Vision and Pattern Recognition*, 2016.
- [65] P. Pérez, M. Gangnet, and A. Blake. Poisson image editing. *ACM Transactions on graphics (TOG)*, 22(3):313–318, 2003.
- [66] J. Pont-Tuset, F. Perazzi, S. Caelles, P. Arbeláez, A. Sorkine-Hornung, and L. Van Gool. The 2017 davis challenge on video object segmentation. *arXiv:1704.00675*, 2017.
- [67] S. A. Ramakanth and R. V. Babu. Featurematch: A general annf estimation technique and its applications. *IEEE Transactions on Image Processing*, 23(5):2193–2205, 2014.
- [68] S. A. Ramakanth and R. V. Babu. Seamseg: Video object segmentation using patch seams. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 376–383, 2014.
- [69] A. Roshan Zamir, A. Dehghan, and M. Shah. Gmcp-tracker: Global multi-object tracking using generalized minimum clique graphs. *Computer Vision–ECCV 2012*, pages 343–356, 2012.
- [70] J. Sánchez. Comparison of Motion Smoothing Strategies for Video Stabilization using Parametric Models. *Image Processing On Line*, 7:309–346, nov 2017.
- [71] G. Seguin, P. Bojanowski, R. Lajugie, and I. Laptev. Instance-level video segmentation from object tracks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3678–3687, 2016.
- [72] T. Shiratori, Y. Matsushita, X. Tang, and S. B. Kang. Video completion by motion field transfer. In *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, volume 1, pages 411–418. IEEE, 2006.
- [73] N. C. Tang, C.-T. Hsu, C.-W. Su, T. K. Shih, H.-Y. M. Liao, et al. Video inpainting on digitized vintage films via maintaining spatiotemporal continuity. *IEEE Trans. Multimedia*, 13(4):602–614, 2011.
- [74] P. Tokmakov, K. Alahari, and C. Schmid. Learning motion patterns in videos. *arXiv preprint arXiv:1612.07217*, 2016.
- [75] Y.-H. Tsai, M.-H. Yang, and M. J. Black. Video segmentation via object flow. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3899–3908, 2016.
- [76] H. V. Vo, N. Q. K. Duong, and P. Pérez. Structural inpainting. In *Proceedings of the 26th ACM International Conference on Multimedia*, MM '18, pages 1948–1956. ACM, 2018.
- [77] P. Voigtlaender and B. Leibe. Online adaptation of convolutional neural networks for the 2017 davis challenge on video object segmentation, 2017.
- [78] S. Wang, H. Lu, F. Yang, and M.-H. Yang. Superpixel tracking. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 1323–1330. IEEE, 2011.
- [79] Y. Wexler, E. Shechtman, and M. Irani. Space-time completion of video. *IEEE Transactions on pattern analysis and machine intelligence*, 29(3), 2007.
- [80] S. Xie and Z. Tu. Holistically-nested edge detection. *International Journal of Computer Vision*, pages 1–16, 2017.
- [81] B. Xu, S. Pathak, H. Fujii, A. Yamashita, and



- H. Asama. Spatio-temporal video completion in spherical image sequences. *IEEE Robotics and Automation Letters*, 2(4):2032–2039, 2017.
- [82] N. Xu, B. Price, S. Cohen, J. Yang, and T. Huang. Deep grabcut for object selection. *arXiv preprint arXiv:1707.00243*, 2017.
- [83] N. Xu, L. Yang, Y. Fan, D. Yue, Y. Liang, J. Yang, and T. Huang. Youtube-vos: A large-scale video object segmentation benchmark. *arXiv preprint arXiv:1809.03327*, 2018.
- [84] H. Yang, L. Shao, F. Zheng, L. Wang, and Z. Song. Recent advances and trends in visual tracking: A review. *Neurocomputing*, 74(18):3823–3831, 2011.
- [85] L. Yang, Y. Wang, X. Xiong, J. Yang, and A. K. Katsaggelos. Efficient video object segmentation via network modulation. *algorithms*, 29:15, 2018.
- [86] M. Y. Yang, M. Reso, J. Tang, W. Liao, and B. Rosenhahn. Temporally object-based video co-segmentation. In *International Symposium on Visual Computing*, pages 198–209. Springer, 2015.
- [87] Y. Yang, G. Sundaramoorthi, and S. Soatto. Self-occlusions and disocclusions in causal video object segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4408–4416, 2015.
- [88] S. You, R. T. Tan, R. Kawakami, and K. Ikeuchi. Robust and fast motion estimation for video completion. In *MVA*, pages 181–184, 2013.
- [89] F.-L. Zhang, X. Wu, H.-T. Zhang, J. Wang, and S.-M. Hu. Robust background identification for dynamic video editing. *ACM Transactions on Graphics (TOG)*, 35(6):197, 2016.



**Thuc Trinh LE** is a Ph.D. candidate in Computer Science and Applied Mathematics at the LTCI Lab of Telecom ParisTech, Paris-Saclay University, France. His research is devoted to the development of machine learning techniques to address some

advanced problems in video editing, video segmentation, and video reconstruction.



**Andrés ALMANSA** is a CNRS Research Director at Université Paris Descartes (France) since 2016. He received his MSc and Ph.D. degrees from ENS Cachan (1999, 2002), his MSc and Engineering degrees from Universidad de la República (1995, 1998). He has been previously working with Telecom ParisTech, ENS Cachan (France), Universitat Pompeu Fabra (Spain) and Universidad de la República (Uruguay). His current research interests include image restoration and analysis, subpixel stereovision and applications to earth observation, high quality digital photography and film editing and restoration.



**Yann GOUSSEAU** received the engineering degree from the cole Centrale de Paris, France, in 1995, and the Ph.D. degree in applied mathematics from the University of Paris-Dauphine in 2000. He is currently a professor at Telecom ParisTech. His research interests include the mathematical modeling of natural images and textures, stochastic geometry, computational photography, computer vision, image, and video processing.



**Simon MASNOU** is a Professor in Mathematics at Claude-Bernard Lyon 1 University (France) and Head of Camille Jordan Institute. His research interests include image processing, shape optimization, calculus of variations, and geometric measure theory.