

# FUN-SIS: a Fully UNsupervised approach for Surgical Instrument Segmentation

Luca Sestini<sup>a,b,\*</sup>, Benoit Rosa<sup>a</sup>, Elena De Momi<sup>b</sup>, Giancarlo Ferrigno<sup>b</sup>, Nicolas Padoy<sup>a,c</sup>

<sup>a</sup>*ICube, University of Strasbourg, CNRS, IHU Strasbourg, France*

<sup>b</sup>*Department of Electronics, Information and Bioengineering, Politecnico di Milano, Milano, Italy*

<sup>c</sup>*IHU Strasbourg, Strasbourg, France*

Automatic surgical instrument segmentation of endoscopic images is a crucial building block of many computer-assistance applications for minimally invasive surgery. So far, state-of-the-art approaches completely rely on the availability of a ground-truth supervision signal, obtained via manual annotation, thus expensive to collect at large scale. In this paper, we present FUN-SIS, a Fully-UNsupervised approach for binary Surgical Instrument Segmentation. FUN-SIS trains a per-frame segmentation model on completely unlabelled endoscopic videos, by solely relying on implicit motion information and instrument *shape-priors*. We define *shape-priors* as realistic segmentation masks of the instruments, not necessarily coming from the same dataset/domain as the videos. The *shape-priors* can be collected in various and convenient ways, such as *recycling* existing annotations from other datasets. We leverage them as part of a novel generative-adversarial approach, allowing to perform unsupervised instrument segmentation of optical-flow images during training. We then use the obtained instrument masks as pseudo-labels in order to train a per-frame segmentation model; to this aim, we develop a *learning-from-noisy-labels* architecture, designed to extract a clean supervision signal from these pseudo-labels, leveraging their peculiar noise properties. We validate the proposed contributions on three surgical datasets, including the MICCAI 2017 EndoVis Robotic Instrument Segmentation Challenge dataset. The obtained fully-unsupervised results for surgical instrument segmentation are almost on par with the ones of fully-supervised state-of-the-art approaches. This suggests the tremendous potential of the proposed method to leverage the great amount of unlabelled data produced in the context of minimally invasive surgery.

## 1. Introduction

Minimally Invasive Surgery (MIS) has established itself as an advantageous alternative to standard open-surgery in several surgical specialties, such as pancreatic and hepatic resections (Chen et al. (2018)), cholecystectomy (Antoniou et al. (2014); Coccolini et al. (2015)), appendectomy (Biondi et al. (2016)) and inguinal hernia (Takayama et al. (2020)). Advantages of MIS mainly derive from the small incisions through which procedures are performed, resulting in several benefits for the patients, such as reduced pain, shorter hospitalization time and less risks of infection. However, together with its benefits, MIS has also introduced new challenges for the surgeons, such as a significantly reduced field-of-view and a complex hand-eye coordination, contributing to an overall increased cognitive workload and a prolonged learning curve (Harrysson et al. (2014)). In order to tackle these challenges, Computer-Assistance has strongly developed in recent years, with the aim to support surgeons through a broad spectrum of applications, including automatic surgical skill analysis (Zia and Essa (2018)), surgical phases segmentation (Twinanda et al. (2016)), tool-tissue interaction estimation (Nwoye et al. (2020)), surgical scene reconstruction (Long et al. (2021)), field-of-view expansion (Bano et al. (2020)), safety checkpoint evaluation (Mascagni et al. (2021)). For

most of these high-level tasks, a crucial building-block is represented by the precise localization of surgical tools in the image space, mainly by pixel-wise classification (i.e. image segmentation).

State-of-the-art approaches for surgical tool segmentation use Deep Learning in order to learn a direct and general mapping between input frames and segmentation masks, robust to challenging factors such as motion blur, occlusions, cluttered background and varying lighting conditions. However, despite the unprecedented results provided by Deep Learning, the problem is still far from being solved for real-world applications: current state-of-the-art Deep Learning approaches rely heavily on manual annotations, which are expensive to obtain at a scale large-enough to allow generalization to real-world scenarios.

Alternatives to standard *in-house annotate & train* pipelines have been proposed, trying to address the annotation problem by cutting the cost of labels, for example by acquiring them through crowd-sourcing platforms (Maier-Hein et al. (2016)) or by generating semi-synthetic datasets with automatically obtained labels (Garcia-Peraza-Herrera et al. (2021)). General object segmentation has been tackled in an unsupervised way when video data are available, such as in Video Object Segmentation (VOS), mainly by leveraging the hypothesis of incoherent background motion, uncorrelated with the foreground (Yang et al. (2019a)). However, state-of-the-art VOS approaches, as they strongly rely on such an hypothesis, tend to fail in the surgical scenario, where foreground (surgical

\*Corresponding author:

*e-mail*: sestini@unistra.fr (Luca Sestini)

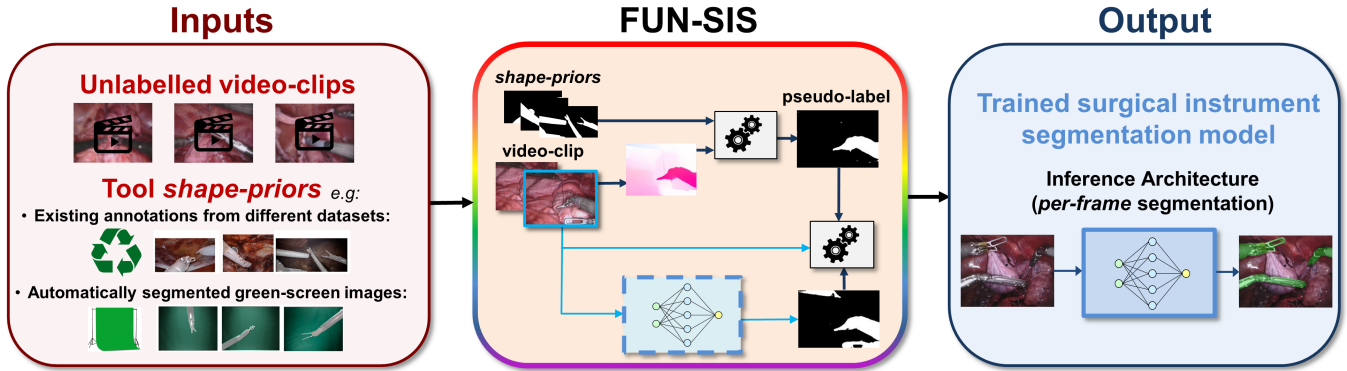


Fig. 1. Paper contribution from the input-output point-of-view. The proposed FUN-SIS approach allows to train a model for surgical tool segmentation requiring as inputs only unlabelled video-clips and tool *shape-priors*, obtainable in various convenient ways (e.g. by recycling existing annotations from other datasets). The method is based on a novel approach for unsupervised surgical tool segmentation of optical-flow images, generating pseudo-label masks, and a newly designed *learning-from-noisy-labels* strategy, allowing to extract a clean supervision signal to train a per-frame segmentation model.

tools) and background (tissue) strongly interact, resulting in coherent and correlated motion.

### 1.1. Contribution

In this paper, we present FUN-SIS, a novel Fully-UNsupervised approach for binary Surgical Instrument Segmentation. The proposed approach allows to effectively train a binary surgical tool segmentation model on completely unlabelled endoscopic videos, solely relying on implicit motion information and instrument *shape-priors*. We define *shape-priors* as binary segmentation masks of the target object, not necessarily coming from the same dataset/domain as the videos. In the specific case of surgical tool segmentation, *shape-priors* can be obtained in convenient and various ways, such as projecting 3D virtual/CAD model of instrument on the image-space, automatically segmenting green-screen recordings, or using existing annotations from other datasets (Figure 1).

In order to achieve this, we make the following contributions:

- we propose a new *generative-adversarial* approach for surgical tool segmentation of optical-flow images, based on simultaneous generation and segmentation of optical-flow images from the *shape-priors*. Compared to state-of-the-art Video Object Segmentation approaches, we relax the hypothesis of incoherent background motion, generally not verified in the surgical domain, letting the *generative-adversarial* training process adapt to the domain characteristics. This leads to state-of-the-art results both on surgical and general Video Object Segmentation datasets;
- we extensively investigate the noise properties of the segmentation masks generated using the proposed optical-flow segmentation approach (*pseudo-labels*), and their impact on neural-network training. We identify and thoroughly analyze two notable properties, namely *unpredictability* and *polarization*, and show that they can be exploited to largely improve segmentation results;

- we propose a novel *learning-from-noisy-labels* strategy, based on an extended *Teacher-Student* approach, allowing to train a *Student* model only on *probably* well-labelled regions of the noisy pseudo-labels. Differently from existing approaches, usually requiring a *Teacher* model trained on clean labels, we carry out an efficient region selection in a fully-unsupervised way, exploiting the aforementioned noise properties. The proposed approach leads to high-quality segmentation results on several surgical datasets, including the popular EndoVis 2017 Instrument Segmentation dataset, while not requiring any ground-truth annotation for the training data.

## 2. Related Work

### 2.1. Surgical Tool Segmentation

Surgical tool segmentation is the task of labelling each pixel of an image as belonging to a specific class among *background* and *tool*, using a single *instrument* class in case of binary segmentation. First attempts to solve this problem used hand-crafted image features, machine learning models (Support Vector Machine, Random Forests) and template matching using tool *shape-priors* (Bouget et al. (2015), Rieke et al. (2016)). Nowadays, research works mostly address the problem of surgical tool segmentation using fully-supervised Deep Learning approaches, which have proved to outperform other existing methods (Bodenstedt et al. (2018)). In particular, Convolutional Neural Network (CNN) architectures have been widely adopted. Garcia-Peraza-Herrera et al. (2017) propose a multi-scale and holistically nested CNN light-weight architecture trained with Dice loss function; Shvets et al. (2018a) modify a VGG16 architecture by adding skip connections, winning the 2017 MICCAI EndoVis Robotic Instrument Segmentation challenge (Allan et al. (2019)). Pakhomov et al. (2019) leverage deep residual learning and dilated convolutions for binary and part tool segmentation. Hasan and Linte (2019) propose a variation of the standard U-Net architecture (Ronneberger et al. (2015)) having a modified decoder and an improved augmentation pipeline. Multi-task learning has also been explored by Laina et al. (2017), by simultaneously

learning segmentation and tool landmarks localization, and by [Islam et al. \(2021\)](#), by introducing the Spatio-Temporal Multi-Task Learning (ST-MTL) model, for surgical instrument segmentation and task-oriented saliency detection. [Jin et al. \(2019a\)](#) propose an attention based approach leveraging motion information, in the form of optical-flow, to improve segmentation accuracy. [Ni et al. \(2020\)](#) propose a bilinear attention network with an adaptive receptive field to tackle the challenges of scale and illumination inter-frames variability. [Kurmann et al. \(2021\)](#) propose an alternative to standard semantic segmentation, first extracting instrument instances and then independently classifying them, reaching state-of-the-art results for this task. Despite the good results obtained by fully-supervised methods, their application is inherently limited by the need for manual annotations, which prevents their scalability. In order to mitigate this problem, [Garcia-Peraza-Herrera et al. \(2021\)](#) produce semi-synthetic samples, merging automatically segmented tools from green-screen recordings and real surgical background images. In the context of robotic surgery, [Colleoni et al. \(2020\)](#) propose the combined use of recorded kinematics and green-screen, in order to cheaply obtain ground-truth segmentation masks. Several works have also tried to tackle the segmentation problem by including synthetic or unlabelled data, in combination with generative approaches. [Sahu et al. \(2020\)](#) propose Endo-Sim2Real, a consistency-based framework for joint training from simulated and unlabelled real data. [Colleoni and Stoyanov \(2021\)](#) propose a cycle Generative Adversarial Network (cycle-GAN) approach to convert synthetic tools into real-looking ones, to be then blended with surgical background images, to form semi-synthetic samples. [Ross et al. \(2018\)](#) pre-train a CNN on unlabelled data, by means of a pretext task carried out using a cycle-GAN architecture, showing a significant boost in segmentation accuracy. [Kalia et al. \(2021\)](#) incorporate unlabelled data in the training process, by mapping annotated frames to the unlabelled data domain using a cycle-GAN architecture, allowing for better generalization to the unlabelled domain. [Marzullo et al. \(2021\)](#) use a *pix2pix* GAN to generate synthetic surgical images from rough segmentation mask of surgical instruments and tissues. In the context of robotic surgery, [Pakhomov et al. \(2020\)](#) record synchronized surgical videos and kinematic joint values and then use the latter to generate synthetic annotations, projecting the estimated tool 3D shapes, obtained via forward kinematics, onto the image space; in order to take into account the possible inaccuracy of the tool model, the segmentation problem is formulated as unpaired image-to-image translation, using a cycle-GAN architecture. An alternative proposed solution to reduce the need for manual annotations is represented by semi-supervision using label propagation. [Zhao et al. \(2020\)](#) propose a flow prediction and compensation framework for semi-supervised tool segmentation, propagating low hertz annotations to unlabelled data using optical-flow. Finally, an unsupervised approach is proposed by [Liu et al. \(2020a\)](#), which generate tool pseudo-labels using handcrafted cues, such as color distribution, and then refine segmentation results exploiting feature correlation between adjacent video frames.

In this work we propose a fully-unsupervised approach for surgical instrument segmentation. Differently from [Pakhomov et al. \(2020\)](#), we do not make use of synchronized kinematic information, making the approach applicable to non-robotic domains (e.g. manual laparoscopy) and to unlabelled video-only datasets (e.g. EndoVis 2017 dataset). In addition, differently from [Liu et al. \(2020a\)](#), we do not rely on domain-specific handcrafted cues, making the approach more robust, flexible and easy to apply to different surgical domains.

## 2.2. Video Object Segmentation

Motion is an important information which is used by the human visual system for *perceptual grouping*, the process of organizing the visual information in order to efficiently perceive and interact with the world. In the general object segmentation framework, as well as for surgical tool segmentation, motion can be a very discriminative cue, easy to obtain from unlabelled videos by means of the available powerful optical-flow estimators. Given the relevance of motion, the computer-vision community has been constantly exploring the task of Video Object Segmentation (VOS). The two standard approaches to it are semi-supervised VOS and unsupervised VOS. Semi-supervised VOS aims to track a target, specified in the first frame of the sequence in the form of a segmentation mask, across the following frames. Unsupervised VOS, instead, aims to separate a salient foreground object from the background. It is worth noticing that, despite its name, unsupervised VOS has often been tackled in literature by means of fully-supervised training (e.g. [Mahadevan et al. \(2020\)](#)): the *unsupervised* attribute indicates, instead, that this family of methods does not need an initial mask of the object, as in semi-supervised VOS. Among the works which have attempted to tackle the unsupervised VOS problem without a ground-truth supervision signal, [Wang et al. \(2017\)](#) propose a geodesic distance based technique, achieving good accuracy, at the cost of high per-frame computation time; more recently, Deep Learning approaches have been proposed: [Yang et al. \(2019a\)](#) propose an adversarial framework to train a neural-network to predict a binary segmentation mask from a frame and the corresponding optical-flow image; [Yang et al. \(2021\)](#) propose an auto-encoder formulation using iterative binding to predict the segmentation mask from optical-flow only. In the context of surgical VOS, the semi-supervised approach is not applicable, due to the repeated changes of instruments during a procedure, and to their motion in and out of the field of view, which would require a continuous re-identification of the objects to be tracked. To our best knowledge, our work represents the first attempt to perform unsupervised VOS of surgical tools, with no annotated ground-truth for training data. The reason for such lack of approaches may lie on the additional challenges that the surgical environment brings to the VOS problem: foreground (tools) and background (tissue) strongly interact with each other, resulting in correlated motion of the two and coherent background motion, thus violating the hypothesis of several state-of-the-art approaches for unsupervised VOS; in addition, tools are not necessarily subject to continuous motion as objects in

general VOS datasets, and may remain still for long periods of time: methods relying on motion segmentation alone, such as Yang et al. (2021), would then fail to capture the object in those sequences.

In this work we propose a novel unsupervised approach for optical-flow tool segmentation, not requiring ground-truth annotations of the training data. In order to tackle the above mentioned challenges, we relax the hypothesis of incoherent background motion, letting a generative-adversarial training process adapt to the domain characteristics. In addition, we show that the pseudo-labels generated from optical-flow tool segmentation, even if noisy, can still provide an effective supervision signal to train a per-frame tool segmentation model, when used in synergy with an efficient *learning-from-noisy-labels* strategy.

### 2.3. Learning from Noisy Labels

Effectively learning from noisy labels is becoming an essential need of Deep Learning applications. In order to gather the massive amounts of annotations required to train Deep Learning models, researchers have recently been looking for alternatives to standard *in-house* annotation, such as crowd-sourcing (Yang et al. (2018)) or automatic-labelling (Guo et al. (2016)). However, while dramatically cutting down the cost of annotations, these approaches tend to provide noisy labels. In order to tackle the *learning-from-noisy-labels* problem, several approaches have been proposed in literature. Following Song et al. (2020), state-of-the-art approaches can be categorized in four groups. *Robust Architecture* methods involve architectural modifications of standard neural networks during training, for example by adding a noise adaptation layer to model the label transition matrix of a noisy dataset (Chen and Gupta (2015)). *Robust Regularization* approaches involve the use of techniques such as data augmentation, weight decay, dropout, and batch normalization to prevent the overfitting of the corrupted examples. *Robust Loss Design* approaches involve the modifications of standard loss functions to make them *noise tolerant*. Examples include generalized cross entropy (GCE, Zhang and Sabuncu (2018)), symmetric cross entropy (SCE, Wang et al. (2019)) and active passive loss (APL, Ma et al. (2020)). Finally, *Sample Selection* approaches, propose strategies to select well-labelled samples. A popular approach for sample selection is multi-network training: MentorNet (Jiang et al. (2018)) uses a mentor network, pre-trained on clean labels, in order to provide a curriculum for the training of a *Student* network. Coteaching (Han et al. (2018)) selects *probably* well-labelled samples according to a *small-loss trick*, training two neural-networks in a collaborative way. While well theoretically motivated, the effectiveness of the above mentioned methods has been proven mainly in the classification task for simpler datasets than the surgical ones, such as artificially modified versions of benchmark datasets like CIFAR (LeCun (1998)), MNIST (Xiao et al. (2017)) and FASHION MNIST (Krizhevsky et al. (2009)), and, less frequently, in real-world datasets with modest-to-medium amount of noise like ANIMAL 10-N (Song et al. (2019)) ( $\approx 8.0\%$  noise rate), Food 101-N

(Lee et al. (2018)) ( $\approx 18.4\%$  noise rate), WebVision (Li et al. (2017)) ( $\approx 20.0\%$  noise rate) and Clothing 1M (Xiao et al. (2015)) ( $\approx 38.5\%$  noise rate). Segmentation differs from standard classification since pixel-labels come grouped in images. This creates the need to rethink standard methods such as *Sample-Selection*, since discarding full labels may represent a waste of useful information. For this reason, local confidence map estimators have been proposed. In the context of 3D medical image segmentation, Yu et al. (2019) propose an approach for semi-supervised learning, where a segmentation model is first trained on clean labels, and then used to produce (noisy) pseudo-labels from unlabelled data, as well as confidence estimations via Monte-Carlo Dropout sampling, in order to train a *Student* model only on well-labelled regions of the pseudo-labels. Nie et al. (2018) also train a segmentation model on clean labels and, in parallel, a confidence model, implemented as a discriminator, in order to discriminate between predicted masks and ground-truth masks, by outputting pixel-wise scores. Unlabelled data can then be fed to segmentation and confidence models to predict pseudo-labels and local confidence maps, enriching the set of labelled training data only with high confidence regions of such predictions. While these methods have achieved good results, their effectiveness is still influenced by the amount and the quality of the available clean labels.

In this work we tackle the problem of learning binary surgical tool segmentation from noisy pseudo-labels obtained from unsupervised segmentation of optical-flow images. Differently from the above mentioned works, our method does not require any set of clean labels in order to perform local region selection on the pseudo-labels. Instead, it leverages favorable properties of the motion-derived pseudo-labels and the finite capacity of neural-networks. These properties and the proposed method will be described in Section 3.

## 3. Proposed Approach

The FUN-SIS approach (Figure 2) is a 3-step method which carries out unsupervised surgical tool segmentation of optical-flow images (step I) and subsequently trains a per-frame segmentation model on the noisy pseudo-labels generated at step I using a new *learning-from-noisy-labels* strategy (steps II and III). The 3 steps are introduced below and detailed in the next sections:

- i) generative-adversarial training of the optical-flow tool segmentation model (called *Teacher*), carried out by simultaneously learning to generate and segment synthetic optical-flow images from tool *shape-priors* (Section 3.1);
- ii) training of a model (called *Proxy*) for tool segmentation of individual frames, using, as direct supervision, the noisy pseudo-labels generated by the *Teacher* model via optical-flow segmentation; the effectiveness of this step is guaranteed by a property of the noise affecting the pseudo-labels, called *unpredictability* (Section 3.2);
- iii) training of a model (called *Student*) for tool segmentation of individual frames, using, as supervision, only *probably*

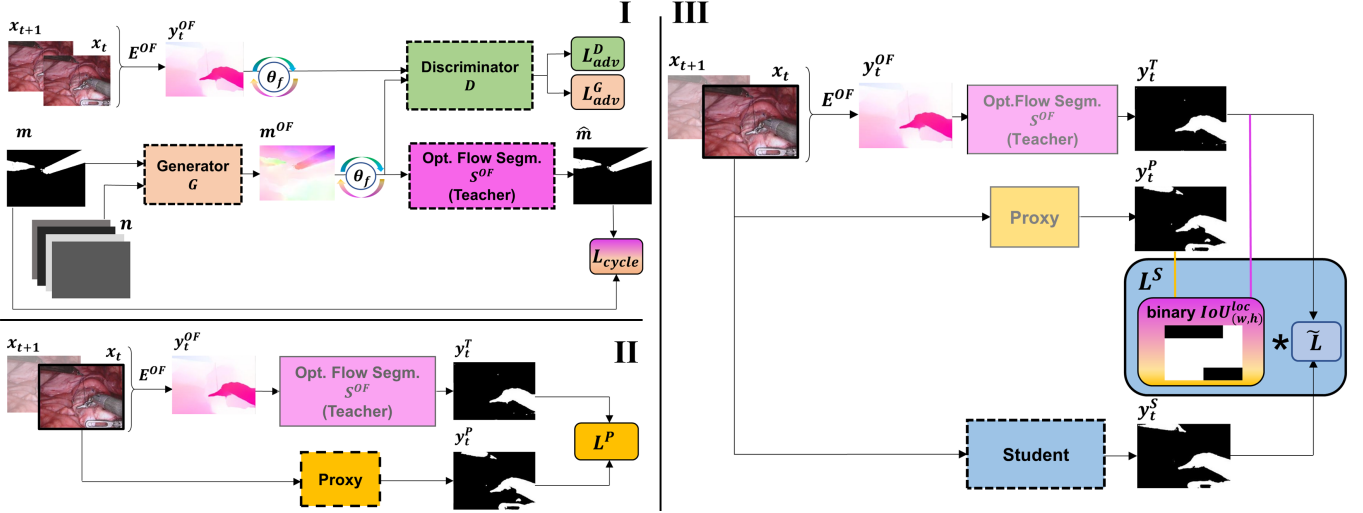


Fig. 2. Overview of proposed FUN-SIS training architecture. **I**: generative-adversarial training of optical-flow segmenter  $S^{OF}$  (Teacher), generator (G) and discriminator (D); generated ( $m^{OF}$ ) and real ( $E^{OF}(x_t, x_{t+1})$ ) optical-flow images undergo augmentation via random rotation  $\theta_f$ ; **II**: Proxy segmentation model training, directly supervised by the pseudo-labels  $y_t^T$ , obtained from optical-flow segmentation by the Teacher model; **III**: Student segmentation model training, leveraging local Intersection-over-Union ( $IoU_{(w,h)}^{loc}$ ) between Teacher and Proxy predictions to select well-labelled regions of  $y_t^T$ .  $\tilde{L}$  is a pixel-wise loss (e.g. cross entropy), masked by the pixel-wise multiplication ( $*$ ) with the binarized local IoU. Loss boxes (L) are color coded to show which models are responsible for their minimization during training. In practice, steps I and II can be carried out simultaneously, as detailed in Section 3.4.

well-labelled regions of the pseudo-labels, selected according to the local agreement between the Teacher and Proxy models; the effectiveness of this step is guaranteed by another property of the noise affecting the pseudo-labels, called *polarization* (Section 3.3).

### 3.1. Step I, Teacher: unsupervised optical-flow segmentation

The proposed approach for unsupervised optical flow-segmentation is based on a generative-adversarial approach, constrained by a cycle-consistency loss. This approach allows to learn the mapping between the domain of optical-flow images and the domain of *shape-priors*, consisting of realistic binary segmentation masks of the target object (in this case surgical tools), without requiring pairwise matching between the two domains. The method is inspired by the classic cycle-GAN architecture (Zhu et al. (2017)), a popular generative architecture for image-to-image translation from unpaired domains. However, it is known that mapping between a domain of minimal complexity, as the binary *shape-priors*, lacking of strong discriminative features, and a more complex one, such as the optical-flow, is an ill-posed problem, suffering from issues such as information-hiding (*‘steganography’* Chu et al. (2017)) and overpowering discriminator, possibly hindering the whole training process.

In order to deal with this *complexity-imbalance*, we propose the following modifications to the standard cycle-GAN:

- we use a single cycle-consistency loss (only for *shape-priors* domain), in order to avoid reconstructing a high-complexity domain sample from a *synthetic* low-complexity domain sample, preventing *‘steganography’*;
- we concatenate the *shape-priors* domain samples with a random noise vector before feeding them to the generator. This allows the generator to produce different

*synthetic* optical-flow images from the same *shape-priors* mask, disentangling the tool silhouette from its motion;

- we make intensive use of on-the-fly image augmentation.

The architecture for the proposed optical-flow segmenter is displayed in Figure 2-I, and discussed below.

Let us consider two consecutive frames belonging to a video,  $x_t, x_{t+1}$  (original frames augmented by an augmentation protocol *AugmData*, consisting of random cropping and flipping), an optical-flow estimator  $E^{OF} : \{x_t, x_{t+1}\} \rightarrow y_t^{OF}$ , where  $y_t^{OF}$  is the optical-flow image in the form of  $[u, v]$  pixel displacement, an optical-flow generator model  $G$ , an optical-flow segmentation model  $S^{OF}$  (also referred to as *Teacher* model, due to its role in steps II and III), a *shape-priors* binary mask  $m$  and a discriminator model  $D$ . The generator  $G$  takes as input the *shape-priors* mask  $m$ , augmented on-the-fly by an augmentation protocol *AugmMask*, consisting of random cropping and flipping, and concatenated with a noise vector  $n$ , sampled from a normal distribution of mean  $\mu$  and standard-deviation  $\sigma$ , and resized to the input mask resolution, and outputs a synthetic optical-flow image  $m^{OF}$ , also in the form of  $[u, v]$  pixel displacement. Both the real and synthetic optical-flow images,  $y_t^{OF}$  and  $m^{OF}$ , undergo on-the-fly augmentation, based on augmentation protocol *AugmFlow*, and following normalization operations:

- **AugmFlow**: the optical-flow is multiplied by a random rotation matrix in the form:

$$R = \begin{bmatrix} \cos \theta_{flow} & -\sin \theta_{flow} \\ \sin \theta_{flow} & \cos \theta_{flow} \end{bmatrix}, \quad (1)$$

where  $\theta_{flow}$  is randomly picked from a uniform distribution. This operation, performed on-the-fly, increases the

variability of the optical-flow, and releases the generator from the burden to generate every possible flow direction;

- **normalization:** each optical-flow image is normalized by dividing it by the maximum pixel displacement  $\sqrt{u^2 + v^2}$  in it. This operation keeps the generated optical-flow image in a controlled range (where maximum displacement has norm equal to 1).

The synthetic optical-flow image  $m^{OF}$  is then fed to the optical-flow segmentation model  $S^{OF}$ , which outputs the *cycled shape-priors* mask  $\hat{m}$ . The real and synthetic optical-flows  $y_i^{OF}$  and  $m^{OF}$  (both augmented and normalized) are fed to the discriminator  $D$ , which is trained to distinguish among the two. Cycle-consistency is ensured by requiring the cycled-mask  $\hat{m}$  to match the input mask  $m$  by means of a standard cross-entropy loss:

$$L_{cycle} = -m \log(\hat{m}) - (1 - m) \log(1 - \hat{m}). \quad (2)$$

Discriminator’s outputs are used to enforce realistic appearance of  $m^{OF}$  by training the discriminator  $D$  and the optical-flow generator  $G$  in an adversarial way. Specifically, the adversarial loss functions are defined as:

$$L_{adv}^G = -\log(D(m^{OF})), \quad (3)$$

$$L_{adv}^D = -\log(1 - D(m^{OF})) - \log(D(y_i^{OF})). \quad (4)$$

The full architecture is trained end-to-end using a standard ADAM optimizer. The discriminator  $D$  is trained to minimize  $L_{adv}^D$ , the optical-flow segmenter  $S^{OF}$  is trained to minimize  $L_{cycle}$ , the optical-flow generator  $G$  is trained to minimize the sum of  $L_{adv}^G$  and  $L_{cycle}$ :

$$L^G = L_{adv}^G + L_{cycle}. \quad (5)$$

### 3.2. Step II, Proxy & the “unpredictability” noise property

The optical-flow segmentation by  $S^{OF}$  (*Teacher* model) is used to generate pseudo-labels for the unlabelled frames: each frame  $x_i$  is paired with the *Teacher*-generated pseudo-label mask  $y_i^T = S^{OF}(y_i^{OF})$ , which is used as direct supervision to train a neural-network (*Proxy* model) to perform tool segmentation of individual frames (Figure 2-II).

The proposed approach to leverage the noisy pseudo-labels relies on findings from Arpit et al. (2017), which show that, while neural-networks are in principle capable of memorizing noisy samples, they tend to first take advantage of shared patterns across training examples, given their finite capacity. In a parallel study, Rolnick et al. (2017) empirically confirmed, in the classification task, that neural-networks can generalize well even when trained on massively noisy data, rather than just memorizing noise, assuming that the noise on a pseudo-label is not conditioned by the corresponding input image itself. We define this condition as the **unpredictability** property.

The noise affecting the pseudo-labels  $y_i^T$  can be divided into two additive processes: the optical-flow estimation noise

and the optical-flow segmentation noise. In both cases, the property of *unpredictability* of noise affecting the pseudo-label  $y_i^T$ , from the single frame  $x_i$ , holds:

- the possible absence of tool motion or presence of background coherent motion in the optical-flow image  $y_i^{OF} = E^{OF}(x_i, x_{i+1})$ , potential sources of  $y_i^T$  noise, cannot be predicted from the individual frame  $x_i$  only, but requires an additional frame ( $x_{i+1}$ ) to be predicted;
- the optical-flow segmentation used to generate the pseudo-labels ( $y_i^T = S^{OF}(y_i^{OF})$ ), second possible source of noise due to the inevitable sub-optimality of  $S^{OF}$  model, does not involve the use of the frame  $x_i$ , contrarily to standard VOS approaches, where both frame and optical-flow are used to make a prediction (e.g. Yang et al. (2019a)).

Given the *unpredictability* property, we can train a neural-network (*Proxy* model) to perform per-frame tool segmentation, using the noisy pseudo-labels  $y_i^T$  directly as supervision signal. The *Proxy* network takes as input the frame  $x_i$  and outputs the segmentation mask  $y_i^P$ . The network is trained to minimize the loss  $L^P$ , which is the sum of the binary cross-entropy loss  $L_{CE}^P$  and the log Intersection-over-Union loss  $L_{IoU}^P$ , weighted by a factor  $\alpha_P$ :

$$L_{CE}^P = -y_i^T \log(y_i^P) - (1 - y_i^T) \log(1 - y_i^P), \quad (6)$$

$$L_{IoU}^P = -\log \frac{\sum(y_i^P y_i^T)}{\sum(y_i^P + y_i^T - y_i^P y_i^T)}, \quad (7)$$

$$L^P = \alpha_P L_{IoU}^P + (1 - \alpha_P) L_{CE}^P. \quad (8)$$

During training, the *Proxy* network, unable to learn the noisy pattern from the pseudo-labels, tries to fit them with the *easiest* compatible pattern, i.e. separating tools from tissue. In order to encourage this effect, we suggest the advantage of using a relatively small-capacity network compared to deeper ones. We experimentally investigate this aspect in our ablation studies, reported in Section 6.2. However, as the training progresses and the pattern is learnt, the loss does not get further minimized, and gradient descent updates remain high, preventing convergence to an optimal solution. This shortcoming is addressed and mitigated at step III below.

### 3.3. Step III, Student & the “polarization” noise property

Together with the *unpredictability* property, a second peculiar property of the noise affecting the pseudo-labels  $y_i^T$  derives from the fact that individual tools, moving coherently, tend to have a uniform appearance in the optical-flow image; this implies that, under ideal conditions (optimal optical-flow estimator  $E^{OF}$ , optimal optical-flow tool segmenter  $S^{OF}$ ), each individual tool will be either perfectly segmented (if moving) or completely mislabelled (if not moving). We define the resulting noise feature as **polarization** property, as a tool can ideally only be perfectly segmented or completely mislabelled by optical-flow segmentation. In the real case,

this property still holds, although occlusions and sub-optimal optical-flow estimation/segmentation tend to inevitably reduce the intensity of the *polarization* (i.e. there will possibly be partially segmented tools). As a practical corollary, the *polarization* property suggests that inside a pseudo-label  $y_i^T$ , there will be either almost-perfectly labelled or almost-completely wrongly-labelled regions. This *polarization* property will be thoroughly investigated in the experiments from Section 6.6. In order to improve training robustness and consistency, we exploit the *polarization* property by designing an unsupervised method to select well-labelled regions of the pseudo-labels  $y_i^T$  (Figure 2-III). The criterion adopted for this selection is the agreement between *Proxy* network predictions  $y_i^P$  (binarized using a threshold value  $\epsilon_P$ ), and pseudo-labels  $y_i^T$  (binarized using a threshold value  $\epsilon_T$ ). The underlying idea is that the *Proxy* network learns a robust general representation (the *easiest* pattern). While its predictions can be incorrect at small-scale (e.g. on border pixels), they are overall reliable at greater scale (i.e. tools are not completely mislabelled as possibly happening in the pseudo-labels). In order to leverage this observation, we introduce a local version of the Intersection-over-Union (IoU) metric, called **local IoU** ( $IoU_{(w,h)}^{loc}$ ). In order to compute  $IoU_{(w,h)}^{loc}$  between two masks, a window of size  $w \times h$  is slid across the masks, using a stride equal to the window size, and IoU is computed inside each time. The output is an image with same resolution as the input masks, whose value at each pixel is the IoU computed for the region containing the pixel (Figure 3). Due to the way it is constructed, it holds that:

$$\frac{1}{W \cdot H} \sum IoU_{(w,h)}^{loc} = IoU, \quad (9)$$

$$\frac{1}{W \cdot H} \sum IoU_{(1,1)}^{loc} = PA, \quad (10)$$

where  $W \times H$  is the size of the input masks, PA is the pixel accuracy metric and the summation is performed over pixels. This makes *local* IoU a metric that interpolates between standard IoU and pixel accuracy, by varying the window size parameter. *Local* IoU is computed between pseudo-label  $y_i^T$  and *Proxy* prediction  $y_i^P$ , and then binarized using a threshold parameter  $\epsilon_{IoU}$ .  $\epsilon_{IoU}$  represents the minimum agreement between *Proxy* and *Teacher* required for a region of  $y_i^T$  to be regarded as well-labelled. The binarized *local* IoU  $\overline{IoU}_{(w,h)}^{loc} = bin(IoU_{(w,h)}^{loc}, \epsilon_{IoU})$  is used to prevent the loss propagation through the *probably* wrongly-labelled regions of the pseudo-labels  $y_i^T$ , during the training of the *Student* network. In particular, the *Student* network takes as input the frame  $x_t$  and outputs the segmentation mask  $y_i^S$ . The network is trained to minimize the loss  $L^S$ , which is the weighted sum of binary cross-entropy loss  $L_{CE}^S$  and log Intersection-over-Union loss  $L_{IoU}^S$ , masked by multiplying each pixel-wise loss by  $\overline{IoU}_{(w,h)}^{loc}$ :

$$L_{CE}^S = \frac{1}{\sum \overline{IoU}_{(w,h)}^{loc}} \overline{IoU}_{(w,h)}^{loc} (-y_i^T \log(y_i^S) - (1 - y_i^T) \log(1 - y_i^S)), \quad (11)$$

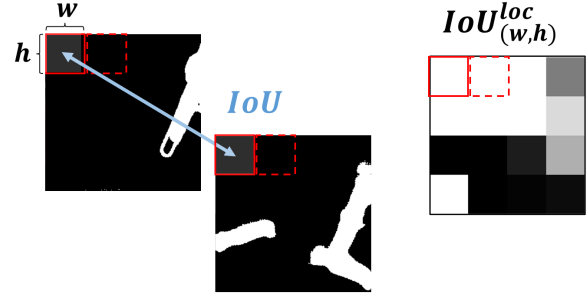


Fig. 3. *Local* IoU  $IoU_{(w,h)}^{loc}$  is computed by sliding a window of size  $w \times h$  on the two input masks, computing standard IoU at each corresponding location. The output is a single-channel image, having the same resolution as the input masks, with each pixel's value being set to the one of the IoU computed for the region it belongs to.

$$L_{IoU}^S = -\frac{1}{\sum \overline{IoU}_{(w,h)}^{loc}} \log \frac{\sum (y_i^S y_i^T \overline{IoU}_{(w,h)}^{loc})}{\sum (y_i^S + y_i^T - y_i^S y_i^T) \overline{IoU}_{(w,h)}^{loc}}, \quad (12)$$

$$L^S = \alpha_S L_{IoU}^S + (1 - \alpha_S) L_{CE}^S. \quad (13)$$

### 3.4. Training Strategy

As presented in Section 3 and shown in Figure 2, the proposed approach involves a 3-step training, where the *Teacher*, *Proxy* and *Student* models are trained successively. However, relying on the hypothesis that a neural-network will not be able to fit the noisy labels, discussed in Section 3.2, we suggest that the *Proxy* network can be trained on the pseudo-labels produced by *Teacher* network while the *Teacher* network is being trained. This allows the training to be a more compact, 2-step process, with steps I and II carried out simultaneously. Comparison between 3-step and 2-step training is reported in Section 5.2.

## 4. Experimental Set-Up

### 4.1. Implementation Details

All models are implemented as neural-networks. Neural-network architectures and hyper-parameters, reported in detail in Appendix A, were determined from preliminary experiments on external data (*phantom* dataset from Sestini et al. (2021)), and kept the same for all the experiments. All the segmentation models have a U-Net-like architecture. The *Proxy* and *Student* networks have slightly different architectures, with the *Proxy* having a 11-convolutional-layer encoder (which we refer to as Unet11) and the *Student* a 16-convolutional-layer (Unet16). Optical-flow estimation was carried out using RAFT (Teed and Deng (2020)), a state-of-the-art approach, trained on the publicly available non-surgical dataset FlyingThings (Mayer et al. (2016)). Training and evaluation were all carried out on  $256 \times 256$  resized versions of the images, regardless of their original resolution/aspect ratio, due to memory constraints. The size of the

noise vector  $n$  was set to 32, and investigated in Section 6.1. Each value of  $n$  was drawn from a normal distribution of mean  $\mu$  equal to 0 and standard-deviation  $\sigma$  equal to 1. The  $IoU_{(w,h)}^{loc}$  window size  $w \times h$  was set to  $64 \times 64$  (1/4 of the image size); the threshold  $\epsilon_{IoU}$  was set to 0.5. An in-depth study regarding  $w$  and  $\epsilon_{IoU}$  was carried out and reported in Section 6.4. The loss balancing factors  $\alpha_P$ ,  $\alpha_S$  from Equations 8&13 were set to 0.8, and investigated in Section 6.3. Augmentations *AugmMask* and *AugmData* were implemented by applying random left-right, up-down flipping and random cropping, with minimal cropped region size equal to  $224 \times 224$ , then bilinearly resampled to  $256 \times 256$ . The angle  $\theta_{flow}$  for the flow rotation in *AugmFlow* was randomly picked in the range  $[-\pi, \pi]$ . All augmentations were applied on-the-fly. Training was carried out using a single NVIDIA Tesla V100 GPU (32 GB). The code will be released upon publication.

#### 4.2. Datasets

In order to validate the proposed contributions, extensive experiments were carried out, both on surgical and general object segmentation datasets. All the data used in our experiments are now presented and categorized as *Video* and *Shape-priors*. Details about their use in the experiments are also reported.

##### Video data:

- **EndoVis2017** (Allan et al. (2019)): dataset from the 2017 MICCAI EndoVis Robotic Instrument Segmentation Challenge. The dataset contains 10 video clips of abdominal porcine procedures, performed using da Vinci Xi systems. Each video contains a total of 300 high-resolution frames ( $1280 \times 1024$ ), recorded at 2 Hz. In the challenge 8x 225 frames were used for training, while the remaining 8x 75 frames and another 2x 300 frames were held out by the organizers for testing. According to the challenge rules, man-made devices not belonging to the da Vinci system (e.g. drop-in Ultra-Sound probe), labelled by the organizers as part of a class called *Other*, are to be included in the *background* class for the binary segmentation task. This introduces the need for a model to perform a semantic differentiation inside the *instrument* class (da Vinci instruments and *Other* instruments), which goes beyond the scope of motion-based segmenters. For this reason we refer to the dataset labelled according to the challenge rules as **EndoVis2017Challenge**, and also consider a second version of it, called **EndoVis2017VOS**, where both da Vinci and other man-made devices are labelled as *instrument*. For the main experiments, we report results on both. We provide results on this dataset according to 2 modalities: 1) following the same evaluation protocol as Shvets et al. (2018a), by performing 4-fold cross-validation on the 8x 225 released training data (regrouped in 4 splits), and reporting the average metric on the 4 splits, for direct and fair comparison with other state-of-the-art approaches; 2) by training on RandSurg, a dataset of

unlabelled data, described below, and testing on the 8x 225 EndoVis2017VOS frames.

- **RandSurg**: this dataset consists of 4 full unlabelled laparoscopic robotic-assisted procedures downloaded from a public repository ([WorldLaparoscopyHospital](#)): adhesiolysis (1036 frames), inguinal hernia repair (1075 frames), appendectomy (500 frames) and ex-vivo suturing demo (525 frames). A set of experiments was carried out by training our model on this dataset and evaluating the performance on EndoVis2017VOS; in order to simulate a realistic application of the FUN-SIS method, and show its ease-of-use, the videos underwent minimal pre-processing (cropping, no trimming, so possibly including out-of-body scenes).
- **STRAS**: this dataset is obtained from endoscopic sub-mucosal dissection procedures performed through the STRAS robotic system (De Donno et al. (2013)), a robotic system consisting of a robotized endoscope, having two lateral channels for flexible robotic tools. The dataset was built from a 5 days-experiment on porcine models<sup>1</sup> (Zorn et al. (2017)), recorded at 30 fps. Each frame was paired with another 1 second apart in the future, for optical-flow computation. The whole dataset was resampled regularly, yielding a total of 5644 frames (~1100 per experiment day). For each day, 200 frames, regularly spaced, were manually annotated for evaluation (1000 annotated samples in total). The dataset contains challenging sequences, involving bleeding, smoke, strong tool-tissue interaction and image blurring. We provide results on this dataset by performing 5-fold cross-validation (each fold corresponding to an experiment day), and reporting the average metric on the 5 splits.
- **Cholec80** (Twinanda et al. (2016)): dataset containing 80 unlabelled videos of manual laparoscopic cholecystectomy procedures captured at 25 Hz and resampled at 1 Hz. We provide qualitative results on this dataset by using the standard split (40 videos for training, 40 videos for testing) to show cross-surgery applicability of the proposed FUN-SIS method.
- **DAVIS2016** (Perazzi et al. (2016)): a popular VOS dataset, containing different moving objects (e.g. animals, people, cars). The dataset consists of 50 clips for a total of 3455 1080p frames with pixel-wise annotations. We provide results on this dataset in order to evaluate the proposed optical-flow segmentation approach on non-surgical videos. To this aim, the standard training-test split was used (30 videos for training and 20 for

<sup>1</sup>The study protocol for this experiment was approved by the Institutional Ethical Committee on Animal Experimentation (ICOMETH No.38.2011.01.018). Animals were managed in accordance with French laws for animal use and care as well as with the European Community Council directive no. 2010/63/EU



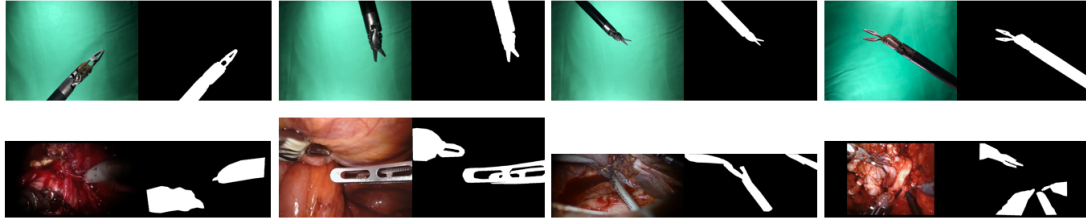


Fig. 4. Examples of *shape-priors* used for the EndoVis2017 experiments, and corresponding source image. Top: tools recorded in front of the green-screen and automatically segmented (Garcia-Peraza-Herrera et al. (2021)), called GrScreenTool; Bottom: frames from multiple robotic-assisted laparoscopic surgeries, manually segmented as part of the RoboTool dataset (Garcia-Peraza-Herrera et al. (2021)). Frames (and also masks) in this dataset come with various resolution/aspect ratios. Note how the appearance of the two domains is different: this is mainly due to the fact that GrScreenTool dataset, recorded using an external camera, show a different point of view on the instruments with respect to the standard surgical camera.

testing), for fair comparison with state-of-the-art VOS approaches.

### Shape-priors:

- **RoboTool:** 514 manually segmented tool masks, from the RoboTool dataset, released by Garcia-Peraza-Herrera et al. (2021). Examples of the original frames and manually segmented tools can be see in Figure 4, bottom. Original masks were cropped to remove the lateral black bands, and resized to  $256 \times 256$  regardless of their original aspect ratio.
- **GrScreenTool:** automatically segmented tools from recordings in front of a green-screen. A total number of 1100 masks were downloaded from the publicly released dataset by Garcia-Peraza-Herrera et al. (2021), mostly having a single tool. Random couples of masks were then selected and merged together, in order to avoid having single-tool masks. Following this strategy, a total number 2200 masks were obtained. Examples of the original green-screen images and extracted tools can be seen in Figure 4, top.
- **STRASmasks:** 2000 projections of approximate 3D virtual/CAD model of the two STRAS tools, used as *shape-priors* in the STRAS experiments; Details regarding the projection operation can be found in Sestini et al. (2021).
- **SegTrackV2** (Li et al. (2013)): 976 manual annotations from the generic VOS dataset SegTrackV2. The dataset includes different segmented objects (e.g. animals, cars, people), used as *shape-priors* in the DAVIS2016 experiments.
- **FBMS59** (Ochs et al. (2013)): 720 manual annotations from the generic VOS dataset FBMS59. The dataset includes different segmented objects (e.g. animals, cars, people), used as *shape-priors* in the DAVIS2016 experiments.

### 4.3. Artificially Corrupted dataset

In order to gain a full understanding of the impact of the noise properties presented in Section 3.2 and 3.3 on

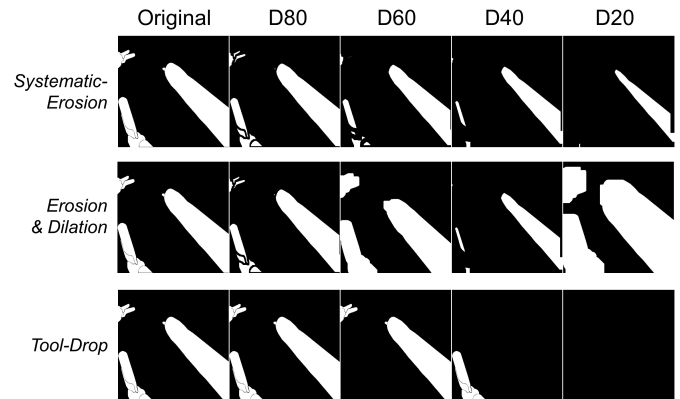


Fig. 5. Samples from the artificially-corrupted versions of EndoVis2017 dataset. From top to bottom: *Systematic Erosion*, *Erosion & Dilation*, *Tool-Drop*. For each noise source, a sample from D80 ( $\sim 80\%$  mean IoU between training sample labels and original ones), D60 ( $\sim 60\%$  mean IoU), D40 ( $\sim 40\%$  mean IoU), D20 ( $\sim 20\%$  mean IoU) is shown.

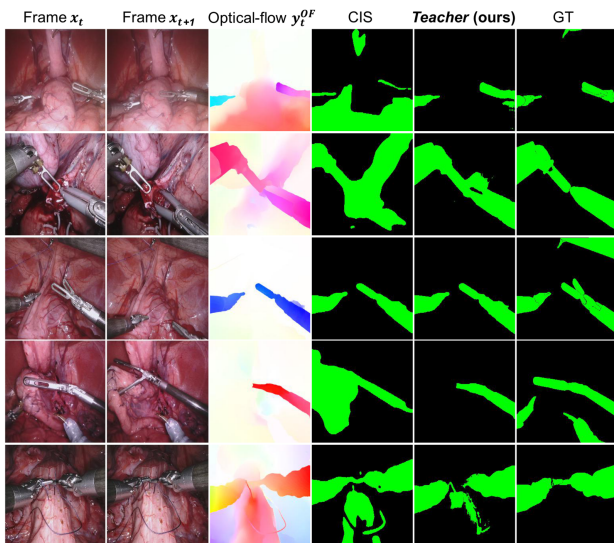
the proposed *learning-from-noisy-labels* approach, we also perform experiments on the EndoVis2017VOS dataset under controlled noise conditions. For these experiments we substitute, in our training pipeline, the pseudo-labels  $y_i^T$  generated by the *Teacher* network, with artificially corrupted versions of the clean labels. To this aim, we consider three types of label corruption, described below:

- *Systematic-Erosion*: each ground-truth mask is eroded;
- *Erosion & Dilation*: each ground-truth mask is randomly eroded or dilated;
- *Tool-Drop*: *full* tool annotations are randomly dropped (i.e. each tool is either *perfectly*-annotated or not-annotated at all).

For each noise type we apply the corresponding transformation, modulating its intensity in order to obtain 4 datasets, {D80, D60, D40, D20}, each one having a mean IoU between the corrupted labels and the ground-truth of  $\sim 80\%$ ,  $\sim 60\%$ ,  $\sim 40\%$ ,  $\sim 20\%$ , respectively (e.g. greater erosion is applied to generate D20 compared to D40, in the *Systematic-Erosion* experiment). Examples of the datasets are shown in Figure 5. We use this dataset as part of the ablation study detailed in Section 6.6.

	Annot. [%]	EndoVis2017	DAVIS2016
Baseline <sub>FS</sub>	100	60.47	73.58
CIS (Yang et al. (2019a))	0*	24.15	60.89 (71.5)
Teacher <sub>RoboTool</sub> (ours)	0	40.08	/
Teacher <sub>GrScreenTool</sub> (ours)	0	<b>40.47</b>	/
Teacher <sub>FBMS</sub> (ours)	0	/	62.72
Teacher <sub>SegTrackV2</sub> (ours)	0	/	<b>63.40</b>

**Table 1. Optical-flow segmentation. Comparison of the proposed method (*Teacher*), using different *shape-priors* for training (RoboTool, GrScreenTool for EndoVis2017VOS experiments; FBMS, SegTrackV2 for DAVIS2016 experiments), with the state-of-the-art CIS approach (without and with post-processing, in parenthesis) and a fully-supervised baseline (Baseline<sub>FS</sub>). Mean IoU [%] is reported. Percentage of annotated training samples required by each method is also reported (Annot. [%]). Note that CIS uses frames and optical-flow to make predictions, while our approach only uses optical-flow.**



**Fig. 6. Optical-flow segmentation on EndoVis2017VOS. Qualitative results showing frame couples used for optical-flow computation, optical-flow images after HSV standard conversion, predictions from CIS (Yang et al. (2019a)) and *Teacher* (trained using RoboTool *shape-priors*), and ground-truth (GT).**

## 5. Experiments and Results Analysis

In this section we present experimental results and comparisons with state-of-the-art methods. First, we analyze the effectiveness of the proposed optical-flow segmentation approach, both on surgical and general object-segmentation datasets. We then analyze results of surgical tool segmentation of individual frames. In order to evaluate model performance, mean Intersection-over-Union (IoU) between predictions and manually annotated ground-truth (GT) is used.

### 5.1. Optical-Flow Segmentation

Optical-flow segmentation by the *Teacher* network was evaluated on EndoVis2017VOS and DAVIS2016, and compared with a state-of-the-art Deep Learning approach for unsupervised Video Object Segmentation, called Contextual Information Separation (CIS, Yang et al. (2019a)), adopting the same evaluation protocol on DAVIS2016. We report the CIS results both with and without post-processing, for fair comparison with our approach which does not make use of it, using the trained network parameters provided by the authors

for DAVIS2016 experiments. Despite being trained using the PWC-net optical-flow estimator (Sun et al. (2018)), we observed that the CIS model provided more accurate results using RAFT-generated optical-flow images: we thus reported results using the latter. On EndoVis2017VOS, the CIS model was trained from scratch, using the RAFT optical-flow estimator: training was carried out using the code publicly released by the authors (Yang et al. (2019b)). We trained our *Teacher* model using RoboTool and GrScreenTool *shape-priors* for EndoVis2017VOS experiment, and SegTrackV2 and FBMS for DAVIS2016 experiment. We also report results of a fully-supervised baseline (Baseline<sub>FS</sub>) model, having the same architecture as the *Teacher* network, trained on GT labels.

Experimental results, presented in Table 1, show that the proposed approach outperforms the state-of-the-art CIS approach (without post-processing) both in the surgical scenario (EndoVis2017VOS dataset) and in general object segmentation (DAVIS2016 dataset). The reason behind the significant improvement on EndoVis2017VOS (+16.32%  $\Delta$ IoU) may reside in the independence of the proposed approach from the hypothesis of incoherent background motion. In fact, our method lets the generator and discriminator adapt to the complexity of the optical-flow domain, generating samples with possible cluttered background and tool occlusion, while still enforcing correct tool segmentation through the cycle-consistency loss. Examples of challenging generated optical-flow images can be seen in Figure 11. Deeper insights on optical-flow generation will be provided the ablation study in Section 6.1. As a result, the optical-flow segmenter becomes more robust to cluttered scenes, where tissue, as well as tools, moves coherently. As shown in the qualitative results shown in Figure 6, the proposed *Teacher* model outperforms the CIS approach especially when tools interact with the anatomy (e.g. pulling tissue, second row from bottom).

### 5.2. Per-frame Surgical Tool segmentation

Per-frame surgical tool segmentation was evaluated on the EndoVis2017Challenge, EndoVis2017VOS and STRAS datasets, according to the modalities reported in Section 4.2, using RoboTool and STRAS Masks as *shape-priors*, respectively. For each experiment, we report results for the following networks:

- *Teacher* network, producing the pseudo-labels  $y_i^T$  from optical-flow segmentation, evaluated against GT masks;

	Annot. [%]	EndoVis2017VOS	EndoVis2017Challenge
TernausNet-16 (Shvets et al. (2018a))	100	(89.06)	83.60 (82.95)
MF-TAPNet (Jin et al. (2019a))	100*	<b>(89.61)</b>	<b>87.56</b> (85.81)
Baseline <sub>FS</sub>	100	88.99	82.55
AGSD (Liu et al. (2020a))	0	(71.47)	67.85 (65.30)
Teacher (ours)	0	40.08	37.03
Proxy (ours)	0	74.78	68.31
<b>Student (ours)</b>	0	<b>83.77</b>	<b>76.25</b>

Table 2. Surgical tool segmentation of individual frames. Comparison of the proposed unsupervised method (trained using RoboTool *shape-priors*), with state-of-the-art unsupervised AGSD approach, fully-supervised approaches TernausNet-16 and MF-TAPNet, and fully-supervised baseline (Baseline<sub>FS</sub>) on the EndoVis2017VOS and EndoVis2017Challenge datasets. Results in parenthesis for state-of-the-art approaches were obtained by training the models using the code released by the authors. Mean IoU [%] is reported. Percentage of annotated training samples required by each method is also reported (Annot. [%]). Note that MF-TAPNet uses 2 consecutive frames at inference time to make a prediction, while the other approaches use individual frames.

	<i>p</i> -value (t-test)	<i>Cohen's d</i>
Proxy-Teacher	$p \ll 0.001$	1.566
Student-Proxy	$p \ll 0.001$	0.612
Baseline <sub>FS</sub> -Student	$p \ll 0.001$	0.448

Table 3. Statistical analysis of tool segmentation results obtained in EndoVis2017VOS (Table 2). For each pair, t-test was run (*p*-values reported in first column) and *Cohen's d* number was computed.

- *Proxy* network, directly trained on the noisy pseudo-labels, producing segmentation masks  $y_i^P$  from individual frames, evaluated against GT masks;
- *Student* network trained using *local* IoU masking, producing segmentation masks  $y_i^S$  from individual frames, evaluated against GT masks. The *Student* network is the output model of the proposed FUN-SIS approach.

For the EndoVis2017Challenge and EndoVis2017VOS experiments we compare the proposed approach with the unsupervised Anchor Generation and Semantic Diffusion (AGSD) approach (Liu et al. (2020a)), based on handcrafted features, and with the fully-supervised state-of-the-art approaches TernausNet-16 (Shvets et al. (2018a)) and MF-TapNet (Jin et al. (2019a)). Results on EndoVis2017VOS for these approaches were obtained by training the models using the code publicly released by the authors (Liu et al. (2020b); Shvets et al. (2018b); Jin et al. (2019b)). Additionally, we compare our results with Baseline<sub>FS</sub>, a model sharing the same architecture as the *Student* network (Unet16), but trained in a fully-supervised way on the GT labels. We do not provide fully-supervised results on the STRAS dataset, due to the lack of GT training labels. We also do not provide results for the unsupervised AGSD approach, due to the fact that the handcrafted cues selected by the authors are specifically tailored for the EndoVis dataset, yielding poor results on the significantly different STRAS dataset.

Experimental results, reported in Table 2, show that the proposed approach enables to effectively train the *Student* network in a fully-unsupervised way, reaching 83.77% IoU on the EndoVis2017VOS dataset, 12.30% above the unsupervised AGSD approach and only 5.22% below the fully-supervised baseline. As hypothesized, the noise affecting the pseudo-

labels generated by optical-flow segmentation cannot be predicted from the individual frames, thus cannot be learnt by the *Proxy* network. This results in a significant improvement of the *Proxy* network’s predictions compared to pseudo-labels used for its training (+34.70%  $\Delta$ IoU on EndoVis2017VOS). On top of this, the *Student* network significantly improves the segmentation quality, by training only on the *probably* well-labelled regions of the pseudo-labels, selected by means of the *local* IoU between pseudo-labels and *Proxy* predictions: the improvement of the *Student* network, with respect to the *Proxy* network, amounts to +8.99%  $\Delta$ IoU on EndoVis2017VOS. Qualitative results presented in Figure 7 clearly show the dramatic improvement of the *Proxy* network compared to the *Teacher* network, and the refining effect of the *Student* network, producing accurate and sharp segmentation masks. In order to assess statistical significance of the results on the EndoVis2017VOS dataset, pairwise t-tests were run (sample size  $N=1800$ ) between *Proxy* & *Teacher*, *Student* & *Proxy* and *baseline<sub>FS</sub>* & *Student*, all showing statistically significant differences ( $p \ll 0.001$  for all the three pairs). In addition, *Cohen's d* number was computed for such pairs, in order to quantify the strength of such statistically significant difference. *Cohen's d* numbers analysis, reported in Table 3, shows that the effect-size of such differences is *very large* between *Proxy* & *Teacher* ( $d > 1.2$ ,  $d = 1.566$ ), *medium/high* between *Student* & *Proxy* ( $0.5 < d < 0.8$ ,  $d = 0.612$ ) and *medium/small* between *fully-supervised baseline* & *Student* ( $0.2 < d < 0.5$ ,  $d = 0.448$ ) (according to Cohen (2013); Sawilowsky (2009)). As expected, the performance on the EndoVis2017Challenge dataset, where devices such as the Ultra-Sound probe are considered as part of the *background* class, is lower than the one on EndoVis2017VOS, while still outperforming the unsupervised AGSD approach (+8.40%  $\Delta$ IoU). This is due to the fact that our approach, despite not being trained using specific *shape-priors* of these tools, is still able to generalize and segment them together with the da Vinci ones. Examples of frames containing the drop-in Ultra-Sound probe are shown in Figure 7, first and fourth row from the top. In order for our approach to learn such semantic discrimination between the two *instrument* classes, pure motion information may not be sufficient. The possible extension of FUN-SIS to multi-class segmentation will be discussed in

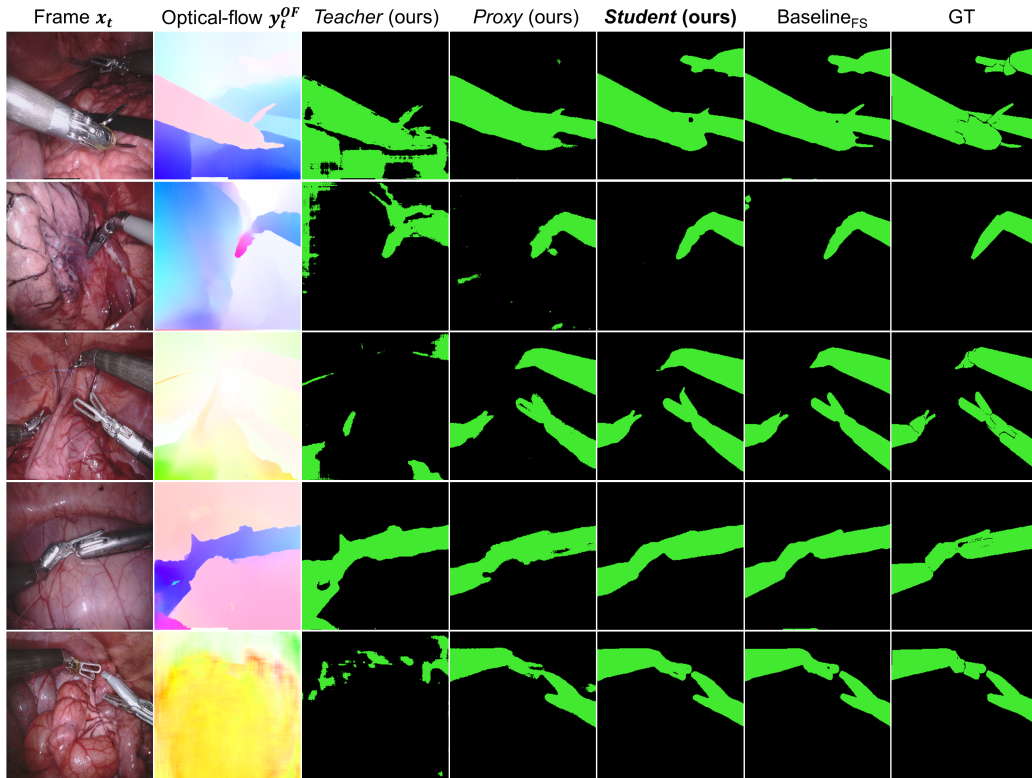


Fig. 7. Surgical tool segmentation on the EndoVis2017VOS dataset. Qualitative results showing, from left to right, input frame  $x_t$ , optical-flow image  $y_t^{OF}$  using HSV standard conversion, predictions from *Teacher* (using RoboTool *shape-priors*), *Proxy*, *Student* and fully-supervised baseline (Baseline<sub>FS</sub>), and ground-truth (GT).

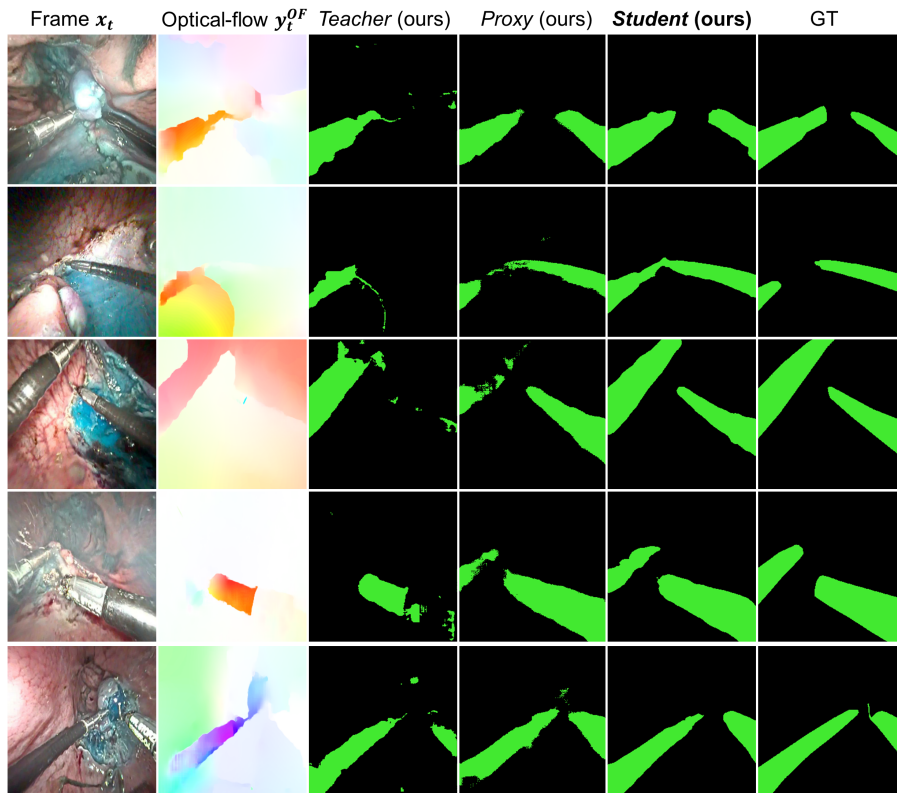


Fig. 8. Surgical tool segmentation on the STRAS dataset. Qualitative results showing, from left to right, input frame  $x_t$ , optical-flow image  $y_t^{OF}$  using HSV standard conversion, predictions from *Teacher* (using STRASMask *shape-priors*), *Proxy* and *Student*, and ground-truth (GT).

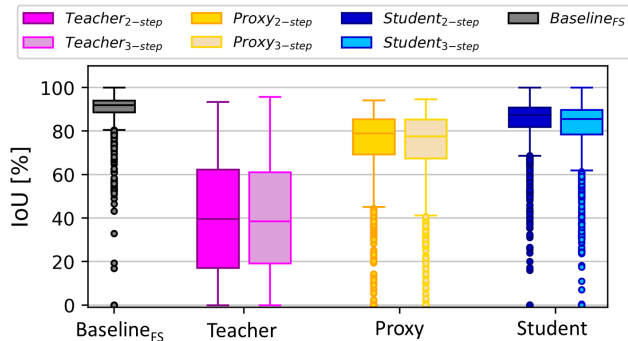


Fig. 9. Box-plots showing IoU distributions from EndoVis2017VOS segmentation experiment (Table 2). Fully-supervised baseline  $Baseline_{FS}$  (grey), *Teacher* (purple 2-step, light purple 3-step), *Proxy* (yellow 2-step, light yellow 3-step), *Student* (blue 2-step, light blue 3-step).

	Annot. [%]	STRAS
Teacher (ours)	0	29.93
Proxy (ours)	0	55.07
<b>Student (ours)</b>	0	<b>66.37</b>

Table 4. Surgical tool segmentation of individual frames. Results of the proposed method on the STRAS dataset using STRAS Masks *shape-priors*. Mean IoU [%] is reported. Percentage of annotated training samples required by each method is also reported (Annot. [%]).

Section 7. We also analyze the difference between the 2-step and 3-step training strategies described in Section 3.4. Results, shown in Figure 9, confirm that the two modalities provide comparable results, as suggested in Section 3.4. We thus consider the 2-step approach superior, due to the shorter training time required. Results obtained on the challenging STRAS dataset, reported in Table 4, confirm the ability of the method to effectively learn surgical tool segmentation in a fully-unsupervised way. The *Student* network, trained without any domain-specific hyper-parameter tuning, reaches an IoU equal to 66.37%, despite being trained on very low-quality pseudo-labels (29.93% IoU). As observable from Figure 8, in fact, optical-flow images appear less sharp compared to the EndoVis2017 ones, mainly due to image blurring and lower image resolution, influencing the overall performance. The implications of the method’s dependency on optical-flow quality will be discussed in Section 7. Additional qualitative results for the *Student* network on the EndoVis2017VOS and STRAS datasets are displayed in Figures 20&21, at the end of the manuscript.

## 6. Ablation Studies and Additional Experiments

In order to provide a more in-depth understanding of the proposed FUN-SIS approach, we performed several ablation studies on crucial aspects of the method.

### 6.1. Optical-Flow Augmentation and Noise Vector Size

We first analyze optical-flow surgical tool segmentation by the *Teacher* network. In particular, we evaluate the impact of the two proposed strategies to tackle the *complexity-imbalance* between optical-flow and *shape-priors* domain in

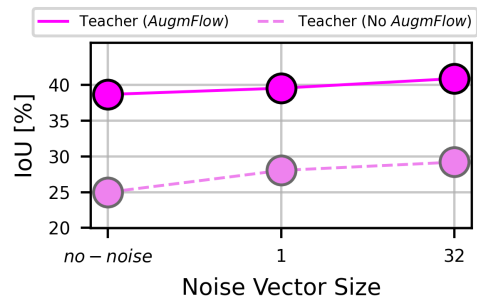


Fig. 10. Analysis of the impact of noise vector size (no-noise, 1, 32) and flow augmentation *AugmFlow* on optical-flow segmentation results by the *Teacher* network on EndoVis2017VOS. Mean IoU [%] is reported.

the generative part of the *Teacher* training, described in Section 3.1: noise concatenation and optical-flow augmentation *AugmFlow*. We trained the *Teacher* model using different sizes of the concatenated noise vector  $n$ , with and without the optical-flow augmentation *AugmFlow*.

Results shown in Figure 10 highlight how optical-flow augmentation *AugmFlow* plays a crucial role in counteracting *complexity-imbalance*, allowing to reach quasi-optimal performance even without noise concatenation (continuous line, “no-noise”). Noise concatenation also appears effective, with peak *Teacher* performance reached with noise size 32 and *AugmFlow*. From qualitative results shown in Figure 11, it can be noted how noise concatenation allows to both generate more realistic and variable optical-flow images and disentangle tools configuration and optical-flow appearance. Note how, when changing *shape-priors*, optical-flow image appearance changes when noise is not concatenated (x0, first block), but remains similar in case of noise concatenation (x1 and x32, second and third block, respectively). It can also be observed how the most variable results are obtained with a noise vector size of 32 (third block, x32), with complexity increasing from leftmost column (noise vector of zeros, more frequently sampled during training) to rightmost column (noise vector of ones, rarely sampled during training), where tools are hardly recognizable.

### 6.2. Proxy Network Architecture

In Section 3.2 we hypothesized the benefit of a limited *Proxy* network capacity, in order to encourage the learning of the *easiest* pattern shared between training samples, compatible with the pseudo-labels. We investigate this hypothesis by evaluating the performance of *Proxy* and *Student* networks, when using different *Proxy* architectures (Unet11 and Unet16, defined in Section 4.1) on the EndoVis2017VOS dataset.

Results shown in Figure 12 confirm that a shallower *Proxy* network (Unet11) learns more effectively from the pseudo-labels than a deeper one (Unet16), quantified in an improvement of +5.44%  $\Delta IoU$ . Additionally, this study provides the experimental proof that the *Student*’s improvements with respect to the *Proxy* are not due to their different architectures.

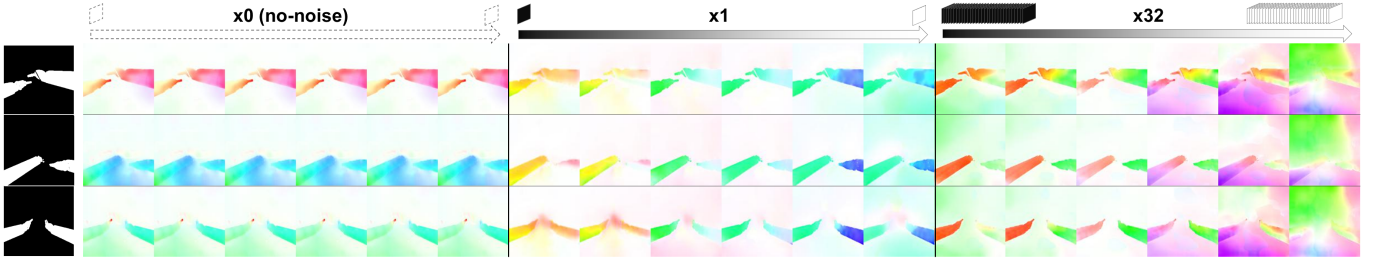


Fig. 11. Qualitative results of the optical-flow generator ( $G$ ), trained using different size of input noise vector among {no-noise,1,32}. First column: input *shape-priors*; first block ( $x_0$ ), no noise concatenation; second block ( $x_1$ ), noise vector of size 1; third block ( $x_{32}$ ), noise vector of size 32. For each of the 3 blocks, from left to right, the noise vector was smoothly interpolated between all zeros to all ones (trivial for  $x_0$ , having no concatenated noise).

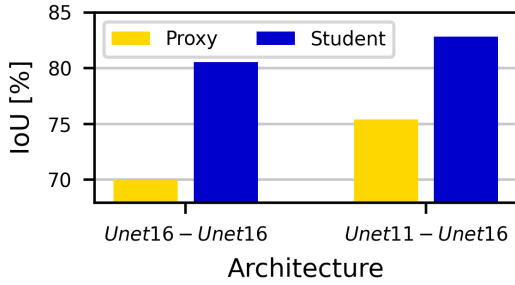


Fig. 12. Analysis of the impact of *Proxy* network’s architecture on surgical tool segmentation results of *Proxy* (yellow) and *Student* (blue) networks, on EndoVis2017VOS. Mean IoU [%] is reported.

Indeed, when using the same architecture for both of them (Figure 12, left) the *Student* still outperforms the *Proxy* by a large margin (+10.61%  $\Delta$ IoU).

### 6.3. Loss Function Coefficients ( $\alpha_P$ , $\alpha_S$ )

We investigate the impact of the balancing factors  $\alpha_P$  and  $\alpha_S$  between cross-entropy (CE) and log IoU losses in *Proxy* and *Student* networks training (Equations 8&13). In our experiments we consider the case  $\alpha_P = \alpha_S = \alpha$ , with  $\alpha$  ranging from 0 (only CE loss) to 1 (only log IoU loss).

Results shown in Figure 13 highlight the positive impact of log IoU loss, especially on the *Proxy* network (+19.79%  $\Delta$ IoU improvement between  $\alpha = 1$  and  $\alpha = 0$ ). This can be in part explained by the diminished-sensitivity of IoU based losses to class-imbalance. However, the greater improvement brought by the log IoU loss to the *Proxy* network, directly trained on raw pseudo-labels, compared to the *Student* network, may suggest that the log IoU loss is more robust to the noise of motion-derived pseudo-labels. Additional in-depth studies are required to investigate this hypothesis.

### 6.4. Local IoU parameters’ impact

While state-of-the-art *learning-from-noisy-labels* approaches usually require a *Teacher* model trained on clean labels in order to identify well-labelled regions of noisy pseudo-labels, we perform this search in a fully-unsupervised way. As detailed in Section 3.3, *probably* well-labelled

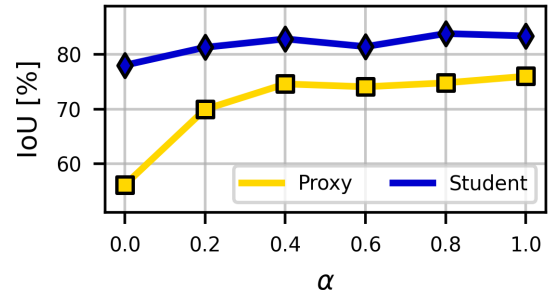


Fig. 13. Analysis of the impact of loss function balancing coefficients  $\alpha_P$  and  $\alpha_S$  on *Proxy* (yellow) and *Student* (blue) networks, on EndoVis2017VOS. We only consider the case  $\alpha_P = \alpha_S = \alpha$ ;  $\alpha$  equal 0 corresponds to cross-entropy loss only,  $\alpha$  equal 1 corresponds to log IoU loss only. Mean IoU [%] is reported.

regions are selected according to the *agreement* between the pseudo-labels (*Teacher* model’s predictions from optical-flow segmentation  $y_t^T$ ) and *Proxy* model’s predictions  $y_t^P$ . The agreement is measured by the *local* IoU, parametrized by the window size  $w$  ( $w = h$  in our experiments), and binarized through the threshold parameter  $\epsilon_{IoU}$ , representing the minimum agreement required to consider a region well-labelled. The choice of these two parameters influences 1) the *effective* number of pixels on which the *Student* network is trained, 2) the average *effective* IoU ( $IoU_{eff}$ ) of the training labels, defined as the IoU between ground-truth masks GT and pseudo-labels  $y_t^T$ , computed only for the selected regions according to the binarized *local* IoU ( $\overline{IoU}_{(w,h)}^{loc}$ ) between  $y_t^T$  and  $y_t^P$ :

$$IoU_{eff} = \frac{|(GT \cap y_t^T) \cap \overline{IoU}_{(w,h)}^{loc}|}{|(GT \cup y_t^T) \cap \overline{IoU}_{(w,h)}^{loc}|}. \quad (14)$$

We evaluate the influence of  $w$  and  $\epsilon_{IoU}$  on the *effective* training size (expressed as total number of selected pixels over total number of pixels in the training dataset) and on the average  $IoU_{eff}$  in the training dataset. For these experiments we considered trained *Teacher* and *Proxy* models on EndoVis2017VOS. We then varied  $\epsilon_{IoU}$  and  $w$  in a grid-like manner, with  $\epsilon_{IoU}$  ranging from 0.0 to 1.0 with a step equal to 0.05, and  $w$  in  $\{1, 2, 4, 8, 16, 32, 64, 128, 256\}$ . For each couple  $(w, \epsilon_{IoU})$  we then evaluated *effective* training size and

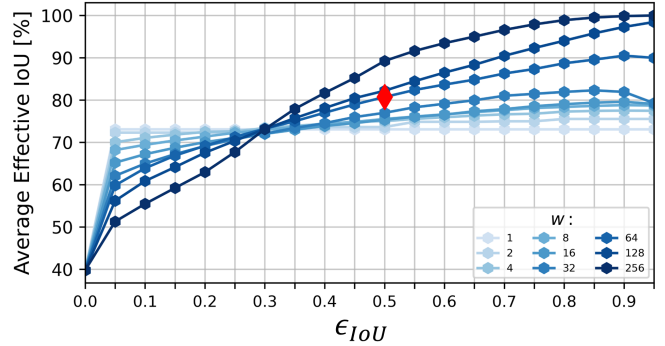
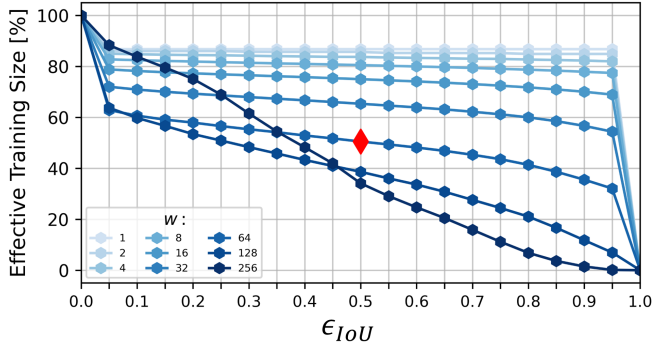


Fig. 14. Impact of local IoU parameters ( $\epsilon_{IoU}$  and window size  $w$ ) on effective training size (left) and average effective IoU (right). x-axis can be interpreted as the level of agreement between *Teacher* and *Proxy* required in order to select a certain region (e.g. with  $\epsilon_{IoU}$  equal to 0.8 a region is considered well-labelled only if the IoU between *Proxy* and *Teacher* predictions for that region is at least 80%). Red markers correspond to  $w = 64$  and  $\epsilon_{IoU} = 0.5$ , the values used in our main experiments.

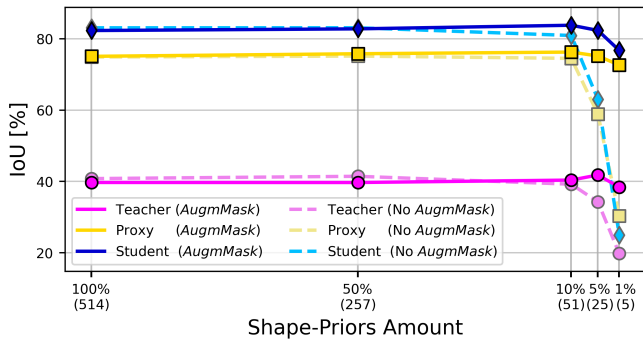


Fig. 15. Analysis of the impact of decreasing *shape-priors* quantity on individual frame and optical-flow segmentation, with and without *AugmMask* augmentation, on EndoVis2017VOS. On the x-axis, the amount of RoboTool *shape-priors* used for training is reported (absolute number and percentage with respect to the total number). Mean IoU [%] for *Student* (blue; dashed: trained without *AugmMask*), *Proxy* (yellow; dashed: trained without *AugmMask*), *Teacher* (purple; dashed: trained without *AugmMask*) is reported.

average  $IoU_{eff}$  on the EndoVis2017VOS training set, in order to provide an insight of the *effective* training carried out.

Experimental results shown in Figure 14 confirm that the agreement between pseudo-labels (optical-flow segmentation masks from the *Teacher*) and *Proxy* predictions is directly correlated to the quality of the pseudo-labels. Figure 14, right, shows the positive correlation between *Proxy-Teacher* agreement ( $\epsilon_{IoU}$ ) and average *effective* IoU, especially for large window sizes  $w$  of the local IoU operation. As expected, the experiment also shows that requiring higher agreement reduces the amount of data effectively used for training, with a similar but inverse relationship. In light of this experiment, the values of window size  $w$  and  $\epsilon_{IoU}$  chosen for experimental validation, respectively 64 and 0.5, represent a good compromise, allowing to train the *Student* network on 50.48% of the total training data on EndoVis2017VOS, with an effective IoU of the pseudo-labels equal to 80.70% (high-quality labels).

Shape-Priors	RoboTool	GrScreenTool
Teacher (ours)	40.08	40.47
Proxy (ours)	74.78	73.63
Student (ours)	83.77	82.63

Table 5. Analysis of the impact of the *shape-priors* dataset on frame segmentation. Comparison of the proposed method trained using RoboTool and GrScreenTool as *shape-priors* on EndoVis2017VOS. Mean IoU [%] is reported.

### 6.5. Shape-Priors Quality & Quantity

*Shape-priors* represent the only external information required by the proposed approach for training. In order to investigate their impact on the whole training process, we performed two sets of experiments. First, we evaluated the performance of our models (*Teacher*, *Proxy*, *Student*) when trained using RoboTool and GrScreenTool *shape-priors*, on EndoVis2017VOS, in order to evaluate the impact of different sources (i.e. *recycled* annotations from a different dataset and automatically segmented tools from green-screen recordings); secondly, we trained our models using different percentages of the available RoboTool *shape-priors*, from 100% to 1%, with and without on-the-fly augmentation *AugmMask*.

Experimental results highlight how our FUN-SIS approach, although it requires *shape-priors* as external source of information, has extremely loose requirements regarding their quality and quantity. Experiments using GrScreenTool, reported in Table 5, provide comparable performance to the ones using RoboTool, despite the significantly different appearance of tools, as shown in Figure 4. In addition, experiments on *shape-priors* quantity (Figure 15), show how the performance of *Teacher*, *Proxy* and *Student* remains optimal even when using as few as 51 RoboTool *shape-priors* masks (10% of total) for training. If augmented on-the-fly using the *AugmMask* protocol (random cropping and flipping), RoboTool *shape-priors* can be further reduced to a total number of 5 instances (1% of total), with limited performance drop (-5.57%  $\Delta IoU$  compared to 100% case).

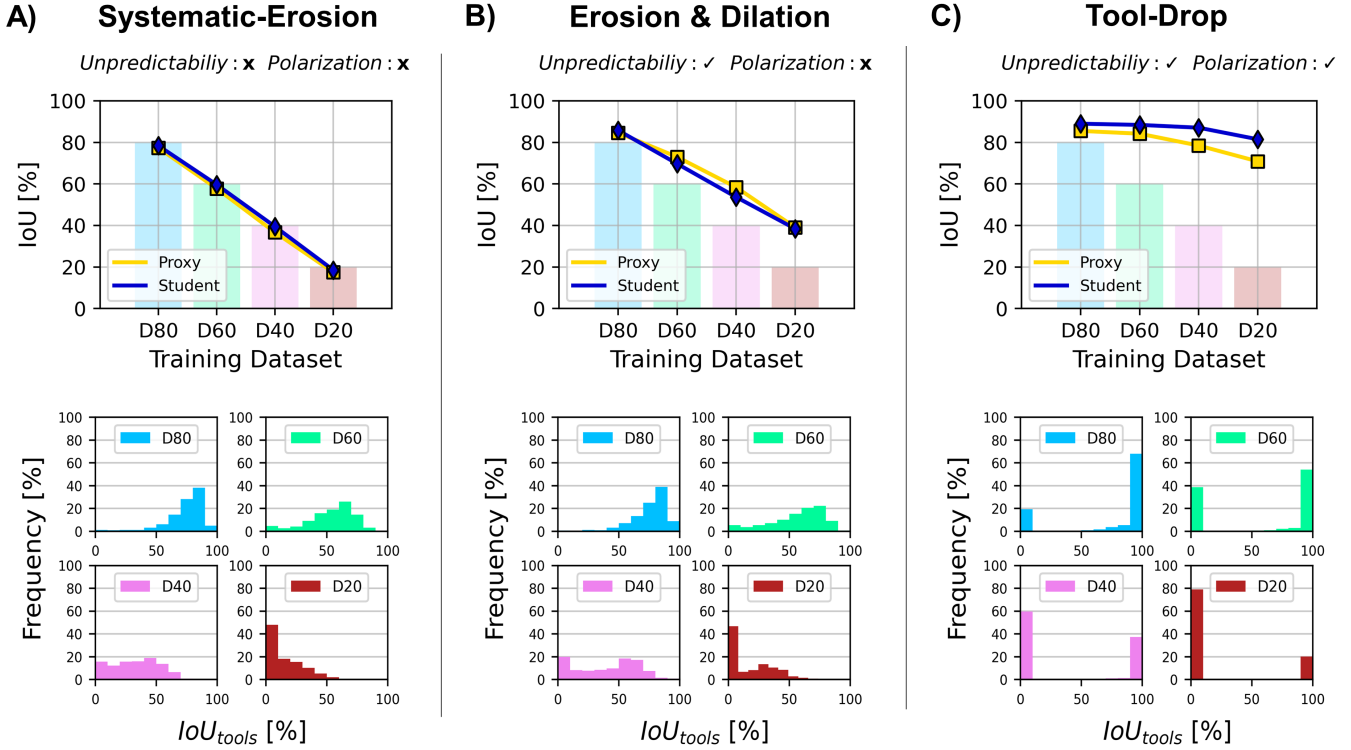


Fig. 16. Analysis of the impact of *unpredictability* and *polarization* noise properties on the proposed method, on the artificially-corrupted EndoVis2017VOS datasets. Top: for each of the 3 noise sources (A, *Systematic-Erosion*, predictable and not-polarized; B, random *Erosion & Dilatation*, unpredictable and not-polarized; C *Tool-Drop*, unpredictable and polarized) *Proxy* (yellow) and *Student* (blue) models were trained on the EndoVis2017VOS training dataset, having ground-truth labels corrupted with different levels of such noise. The colored bars are meant to improve readability, by visually showing the mean IoU between each training dataset labels and ground-truth clean labels ( $\sim 80\%$  for D80,  $\sim 60\%$  for D60,  $\sim 40\%$  for D40,  $\sim 20\%$  for D20); Bottom: for each set of noisy labels, per-tool IoU histograms ( $IoU_{tools}$ ) computed as shown in Figure 17, are reported.

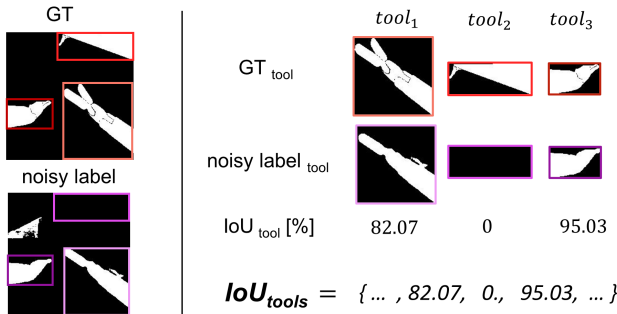


Fig. 17. Computation of per-tool IoU between ground-truth masks and noisy labels. Left: example of ground-truth mask (GT) and noisy label. The smallest region containing each tool in the GT mask is extracted; the same exact region is extracted from the noisy label. Right: Intersection-over-Union ( $IoU_{tool}$ ) is computed between each region extracted from GT ( $GT_{tool}$ ) and noisy label ( $noisy\ label_{tool}$ ); the process is repeated for each tool in each frame of the dataset, and each  $IoU_{tool}$  is stored in  $IoU_{tools}$ . The distribution of per-tool IoU can then be visualized through histogram plots (Figures 16&18).

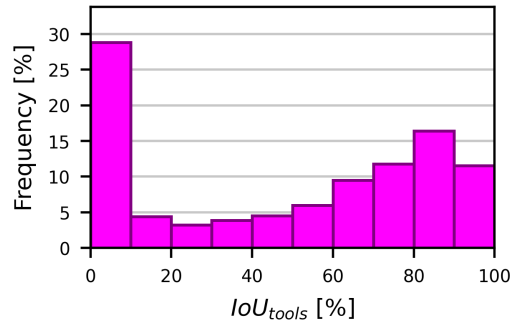


Fig. 18. Per-tool IoU histogram ( $IoU_{tools}$ ), computed as shown in Figure 17, for pseudo-labels derived from motion segmentation by the *Teacher* model on EndoVis2017VOS. Note how the distribution tends to be polarized on leftmost bin (completely mislabelled tools) and rightmost bins (*almost*-perfectly segmented tools).

### 6.6. Noise properties (unpredictability & polarization)

We investigate the impact of the *unpredictability* and *polarization* properties presented in Sections 3.2 and 3.3 on the proposed *learning-from-noisy-labels* approach. To this aim, we carried out experiments with artificially controlled type and intensity of noise affecting the pseudo-labels, as described in Section 4.3. We then substituted the pseudo-labels  $y_i^T$ , in

our training pipeline, with the corrupted EndoVis2017VOS labels and trained the *Proxy* and the *Teacher* networks according to the same modalities as the previous experiments. The three noise strategies presented in Section 4.3 were designed to highlight the effect of the *unpredictability* and *polarization* properties. In *Systematic-Erosion* experiment, each mask was eroded, making the noise signal *predictable* and *not-polarized* (all tools are equally affected by the noise); in *Erosion&Dilatation* experiment, each mask was either randomly



eroded or dilated, making the noise signal *unpredictable*, but still *not-polarized* (each tool mask is affected by an error, either due to erosion or dilation); finally, in *Tool-Drop* experiment, individual tools were either perfectly annotated or not annotated at all, making the noise signal both *unpredictable* and *polarized*.

Results of the conducted experiments (Figure 16) clearly highlight the impact of the two noise properties, as well as the ability of the proposed solution to leverage them. When the noise is predictable (Figure 16-A, top), the *Proxy* network can perfectly learn to fit it, even when the corruption is minimal (D80). Contrarily, when noise cannot be inferred from single frames (Figure 16-B&C, top), the *Proxy* network, unable to learn the noise pattern, will learn the easiest general pattern compatible with the labels, resulting in significantly better predictions than the noisy labels used for its training (on average, +13.76%  $\Delta$ IoU in *Erosion&Dilation*, +29.75%  $\Delta$ IoU in *Tool-Drop*). The effectiveness of the *Student* network training is instead mainly influenced by the *polarization* property. When the noise is not polarized (Figure 16-A&B, top), the *Student* network does not benefit from region selection through *local* IoU (+1.69% and -1.87%  $\Delta$ IoU, respectively, of *Student* compared to *Proxy* network). Instead, when the noise is polarized, well-labelled regions can be effectively identified using *local* IoU, allowing for a consistent improvement of *Student* predictions, compared to *Proxy* ones (+6.73%  $\Delta$ IoU on average, +8.60%  $\Delta$ IoU in D40). The improvement is aligned with the one obtained in the experiments from Section 5.2 (+8.99%  $\Delta$ IoU), where the pseudo-labels were produced via unsupervised surgical tool segmentation by the *Teacher* network and had an IoU with the GT equal to 40.08%. Overall, the proposed approach allows to maintain an IoU of at least 81.49% (compared to the 88.99% reached by fully-supervised training of the *Student* model on clean labels, Table 2), even when trained on extremely low-quality training labels (Figure 16-C, top: *Tool-Drop*, D20 i.e.  $\sim$ 20% IoU between training labels and GT). When trained on D80 and D60, the *Student* network reaches optimal performance (88.98% and 88.41% IoU, respectively).

In order to provide a direct visualization of the *polarization* property, we also report, for each set of noisy labels, including the motion-derived pseudo-labels by the *Teacher* model, per-tool IoU histograms ( $\text{IoU}_{\text{tools}}$ ). Per-tool IoU can be computed, as shown in Figure 17, by extracting the smallest regions containing each tool from the GT labels, and computing the IoU between this region and the corresponding one from the corresponding pseudo-label. This process, while approximate (an extracted region from GT label may contain more than one tool), allows to produce a clear visualization of the *polarization* property, by plotting the histogram of the obtained  $\text{IoU}_{\text{tools}}$ . Histograms are shown in Figure 18, for motion-derived pseudo-labels, and in Figure 16, bottom, for artificially corrupted labels. From Figure 16, bottom, it is possible to intuitively compare the case of not-polarized noise (A,B), where  $\text{IoU}_{\text{tools}}$  values are mostly distributed around a single peak, to polarized noise (C), where the values appear concentrated on leftmost bin (full tool annotations missed) and

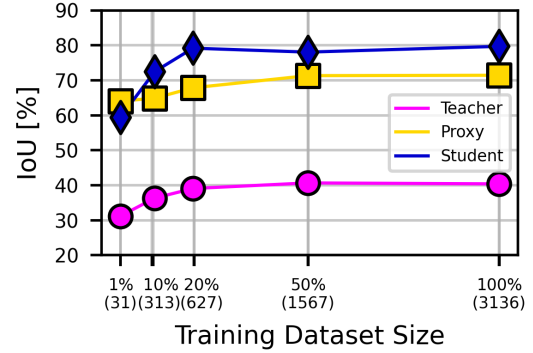


Fig. 19. Analysis of proposed method performance when trained on increasing amounts of unlabelled RandSurg data, a dataset consisting of randomly selected surgical videos, downloaded from the public repository [WorldLaparoscopyHospital](#), and tested on EndoVis2017VOS. On the x-axis, the amount of RandSurg frames used for training is reported (absolute number and percentage with respect to the total number). Mean IoU [%] for *Student* (blue), *Proxy* (yellow), *Teacher* (purple) is reported.

rightmost bin (perfectly labelled tools). In the case of pseudo-labels derived from optical-flow segmentation (Figure 18), the histogram, despite being smoothed by the sub-optimality of optical-flow estimator and segmenter described in Section 3.3, still displays the *polarization* property, allowing efficient *Student* network training.

### 6.7. Random Unlabelled Data

In order to show the ease-of-use and robustness of the proposed FUN-SIS approach, we trained our models on the surgical robotic dataset RandSurg, described in Section 4.2 and tested on EndoVis2017VOS. The RandSurg dataset was created by collecting random public videos of surgical procedures, and performing minimal data curation. Training was carried out according to the same modalities as the other experiments, using RoboTool *shape-priors* and varying amounts of the RandSurg data, ranging from very few (31 i.e. 1% of total available) to all the available frames (3136).

Experimental results shown in Figure 19 show that, despite the limited data curation and pre-processing of the input data, the method can easily leverage the increasing amount of available data to effectively train the models. The *Student* network reaches a peak IoU equal to 79.65% on EndoVis2017VOS, comparable to the 83.77% obtained when training on unlabelled data from the same dataset (Table 2).

### 6.8. FUN-SIS applicability on another domain: Cholec80

We demonstrate the applicability of the proposed FUN-SIS approach on a different domain than the robotic one it was validated on. To this aim, we trained and qualitatively tested our *Student* model on the unlabelled Cholec80 dataset, consisting of manual laparoscopic cholecystectomy procedures. Training was carried out using RoboTool *shape-priors*, despite the different appearance of tools between robotic and manual laparoscopic videos.

Results shown in Figure 22 qualitatively confirm that the proposed method is applicable to a different surgical domain, even without domain-specific hyper-parameters tuning and with minimal pre-processing. Furthermore, they prove that despite the differences between *shape-priors* and target tools, segmentation can still be effectively carried out.

## 7. Discussion and Future Work

In order to validate the proposed FUN-SIS approach, several experiments were performed and presented, including optical-flow segmentation (Section 5.1), per-frame segmentation (Section 5.2, main experiment) and several ablation studies (Section 6), dissecting the method and highlighting its key aspects. The obtained results strongly support the soundness of FUN-SIS: binary surgical tool segmentation was effectively carried out in various datasets including EndoVis2017 (robotic surgery), STRAS (flexible endoscopic surgery), and Cholec80 (manual laparoscopic surgery). When evaluated on EndoVis2017VOS, our *Student* network reaches an IoU of 83.77%, 12.30% above the state-of-the-art unsupervised AGSD approach, and only 5.84% below the state-of-the-art MF-TAPNet approach. Additionally, the proposed unsupervised approach for surgical tool segmentation of optical-flow images outperforms state-of-the-art approaches by a large margin on EndoVis2017VOS (+16.32%  $\Delta$ IoU). Ablation studies proved that the method is extremely robust to the way *shape-priors* are obtained, with no significant performance difference between using automatically segmented tools from green-screen recordings and *recycled* manual annotations from other datasets. In addition, FUN-SIS showed great robustness to limited *shape-priors* quantity, performing optimally on EndoVis2017VOS even using as few as 51 RoboTool *shape-priors* masks for training. Ablation studies highlighted other interesting aspects, as the benefits of using a log Intersection-over-Union loss when training on noisy pseudo-labels, and the effectiveness of the proposed optical-flow augmentation strategy on video object segmentation. Finally, the extensive analysis on pseudo-label noise properties and their impact on neural-network training, as well as the proposed *learning-from-noisy-labels* strategy to leverage them, may serve as base for future work on object segmentation using noisy labels, still largely unexplored. Despite the satisfying results, the proposed work still presents potential room for improvement:

- when selecting well-labelled regions through *local* IoU, a great amount of the available data are currently discarded (49.52% of total available pixels in EndoVis2017VOS experiment). These *uncertainly*-labelled pixels could be exploited with semi-supervised-like strategies, and contribute to the *Student* network training;
- the window used to compute the *local* IoU has fixed dimensions and is slid regularly on the masks with fixed width and stride; a more flexible approach, adapting to the varying tool size and location, may be beneficial;

- the *Proxy* network is subjected to strong gradients while training directly on the noisy pseudo-labels, resulting in possible performance oscillations. This can potentially hinder the *Student* network training, if the *Proxy* network training is stopped in a poor weight parameters configuration. This problem could be mitigated by using approaches such as self-ensembling (Nguyen et al. (2020)), regularizing *Proxy* network training;
- the FUN-SIS performance is overall influenced by the quality of the optical-flow images, which depends, in turn, on the endoscopic camera resolution and the optical-flow estimator. Current research on models for optical-flow computation specifically tailored for endoscopic images, as well as the increasing use of high-definition endoscopic cameras, could naturally contribute to improve the effectiveness of the proposed FUN-SIS method;
- the FUN-SIS approach, completely relying on instrument motion, is currently unable to perform semantic differentiation among the *instrument* class. Strategies to extend FUN-SIS to multi-class segmentation could be explored, possibly involving motion patterns analysis and the use of limited external semantic supervision.

## 8. Conclusion

In this paper we presented FUN-SIS, a novel Fully-UNsupervised approach for Surgical Instruments Segmentation. FUN-SIS effectively trains a per-frame surgical tool segmentation model on completely unlabelled endoscopic videos, solely relying on implicit motion information and instrument *shape-priors*. In order to achieve this, we made several contributions, including a novel unsupervised optical-flow tool segmentation approach and a newly designed *learning-from-noisy-labels* strategy. The proposed contributions were extensively validated on different surgical datasets (flexible endoscopic, robotic and laparoscopic procedures). On the popular MICCAI 2017 EndoVis Robotic Instrument Segmentation Challenge dataset, the proposed unsupervised approach performs almost on par with state-of-the-art fully-supervised models.

In conclusion, we hope that this work can contribute to the development of new segmentation methods requiring reduced supervision for training, fully exploiting the massive amounts of data which minimally invasive surgery can provide.

## Acknowledgments

This work was supported by the ATLAS project. The ATLAS project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 813782. This work was also partially supported by French State Funds managed by the Agence Nationale de la Recherche (ANR) through the Investissements d'Avenir Program under Grant

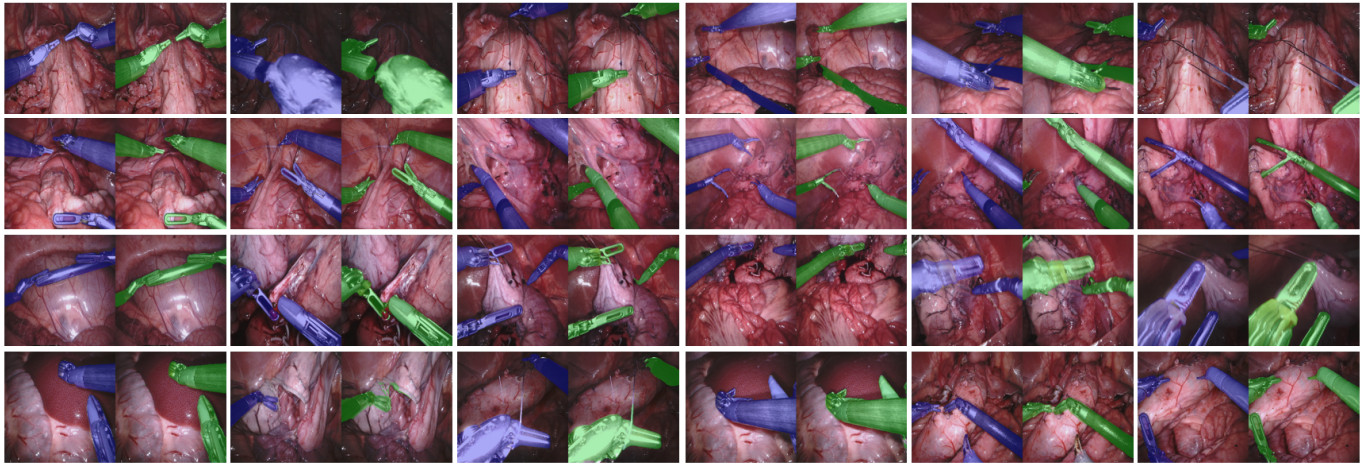


Fig. 20. Qualitative results on the EndoVis2017VOS dataset, from the experiment reported in Table 2. Original frame overlapped with ground-truth (blue) and *Student* network's prediction (green).

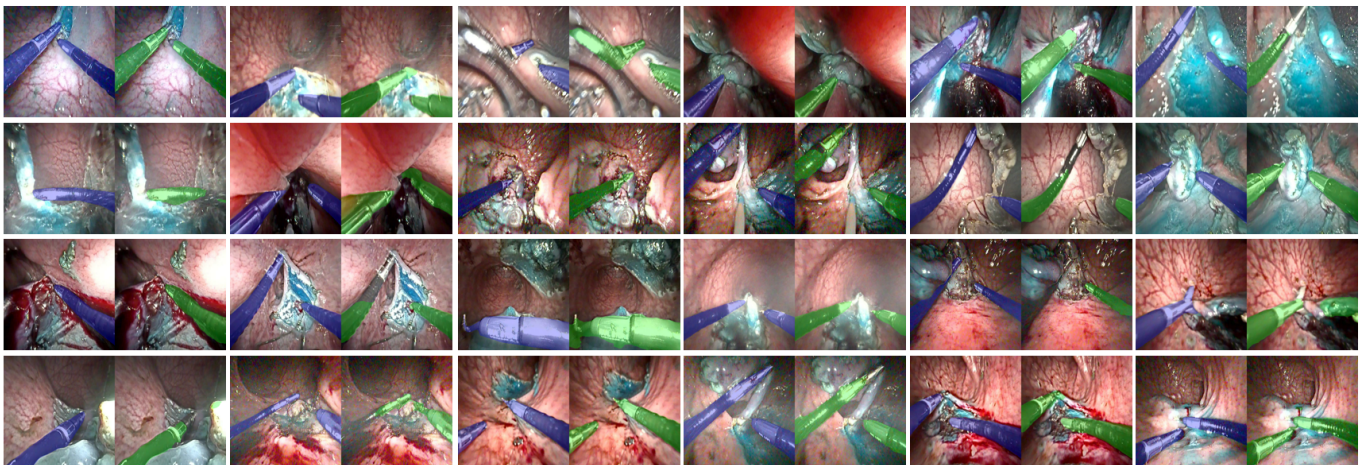


Fig. 21. Qualitative results on the STRAS dataset, from the experiment reported in Table 4. Original frame overlapped with ground-truth (blue) and *Student* network's prediction (green).

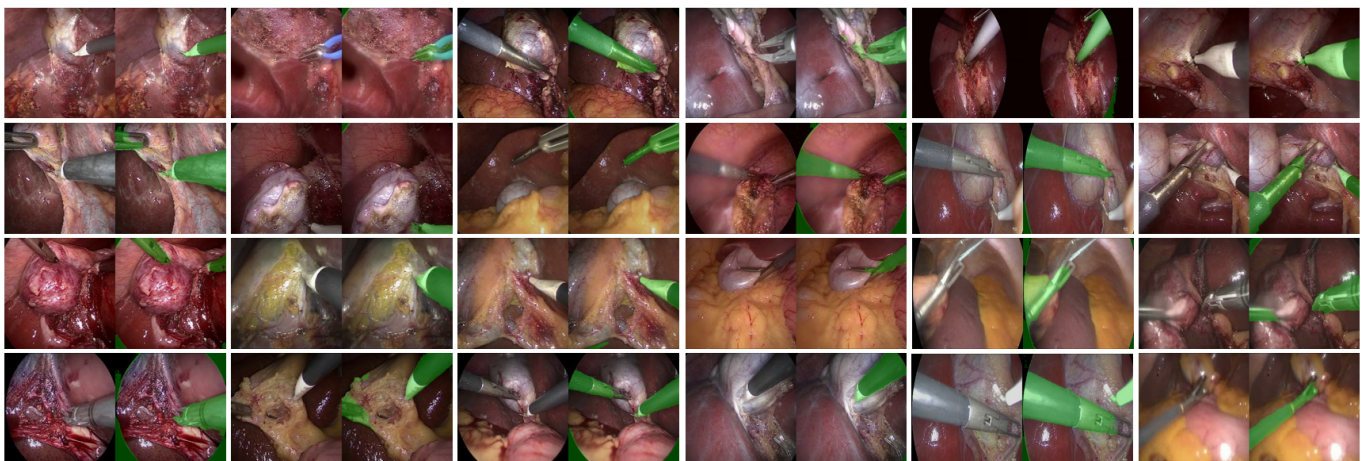


Fig. 22. Qualitative results on the Cholec80 dataset. Original frame and overlapping between *Student* network prediction and original frame are shown. Training was carried out using RoboTool *shape-priors*.

ANR-11-LABX-0004 (Labex CAMI) and Grant ANR-10-IAHU-02 (IHU-Strasbourg).

## References

- Allan, M., Shvets, A., Kurmann, T., Zhang, Z., Duggal, R., Su, Y.H., Rieke, N., Laina, I., Kalavakonda, N., Bodenstedt, S., et al., 2019. 2017 robotic instrument segmentation challenge. arXiv preprint arXiv:1902.06426 .
- Antoniou, S.A., Antoniou, G.A., Koch, O.O., Pointner, R., Granderath, F.A., 2014. Meta-analysis of laparoscopic vs open cholecystectomy in elderly patients. *World Journal of Gastroenterology*: WJG 20, 17626.
- Arpit, D., Jastrzkebski, S., Ballas, N., Krueger, D., Bengio, E., Kanwal, M.S., Maharaj, T., Fischer, A., Courville, A., Bengio, Y., et al., 2017. A closer look at memorization in deep networks, in: *International Conference on Machine Learning*, PMLR. pp. 233–242.
- Bano, S., Vasconcelos, F., Tella-Amo, M., Dwyer, G., Gruijthuijsen, C., Vander Poorten, E., Vercauteren, T., Ourselin, S., Deprest, J., Stoyanov, D., 2020. Deep learning-based fetoscopic mosaicking for field-of-view expansion. *International journal of computer assisted radiology and surgery* 15, 1807–1816.
- Biondi, A., Di Stefano, C., Ferrara, F., Bellia, A., Vacante, M., Piazza, L., 2016. Laparoscopic versus open appendectomy: a retrospective cohort study assessing outcomes and cost-effectiveness. *World Journal of Emergency Surgery* 11, 1–6.
- Bodenstedt, S., Allan, M., Agustinos, A., Du, X., Garcia-Peraza-Herrera, L., Kennigott, H., Kurmann, T., Müller-Stich, B., Ourselin, S., Pakhomov, D., et al., 2018. Comparative evaluation of instrument segmentation and tracking methods in minimally invasive surgery. arXiv preprint arXiv:1805.02475 .
- Bouget, D., Benenson, R., Omran, M., Riffaud, L., Schiele, B., Jannin, P., 2015. Detecting surgical tools by modelling local appearance and global shape. *IEEE transactions on medical imaging* 34, 2603–2617.
- Chen, Q., Merath, K., Bagante, F., Akgul, O., Dillhoff, M., Cloyd, J., Pawlik, T.M., 2018. A comparison of open and minimally invasive surgery for hepatic and pancreatic resections among the medicare population. *Journal of Gastrointestinal Surgery* 22, 2088–2096.
- Chen, X., Gupta, A., 2015. Webly supervised learning of convolutional networks, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1431–1439.
- Chu, C., Zhmoginov, A., Sandler, M., 2017. Cyclegan, a master of steganography. arXiv preprint arXiv:1712.02950 .
- Coccolini, F., Catena, F., Pisano, M., Gheza, F., Fagioli, S., Di Saverio, S., Leandro, G., Montori, G., Ceresoli, M., Corbella, D., et al., 2015. Open versus laparoscopic cholecystectomy in acute cholecystitis. systematic review and meta-analysis. *International journal of surgery* 18, 196–204.
- Cohen, J., 2013. *Statistical power analysis for the behavioral sciences*. Academic press.
- Colleoni, E., Edwards, P., Stoyanov, D., 2020. Synthetic and real inputs for tool segmentation in robotic surgery, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 700–710.
- Colleoni, E., Stoyanov, D., 2021. Robotic instrument segmentation with image-to-image translation. *IEEE Robotics and Automation Letters* 6, 935–942.
- De Donno, A., Zorn, L., Zanne, P., Nageotte, F., de Mathelin, M., 2013. Introducing stras: A new flexible robotic system for minimally invasive surgery, in: *2013 IEEE International Conference on Robotics and Automation*, IEEE. pp. 1213–1220.
- Garcia-Peraza-Herrera, L.C., Fidon, L., D’Ettorre, C., Stoyanov, D., Vercauteren, T., Ourselin, S., 2021. Image compositing for segmentation of surgical tools without manual annotations. *IEEE Transactions on Medical Imaging* 40, 1450–1460.
- Garcia-Peraza-Herrera, L.C., Li, W., Fidon, L., Gruijthuijsen, C., Devreker, A., Attilakos, G., Deprest, J., Vander Poorten, E., Stoyanov, D., Vercauteren, T., et al., 2017. Toolnet: holistically-nested real-time segmentation of robotic surgical tools, in: *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, IEEE. pp. 5717–5722.
- Guo, Y., Zhang, L., Hu, Y., He, X., Gao, J., 2016. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition, in: *European conference on computer vision*, Springer. pp. 87–102.
- Han, B., Yao, Q., Yu, X., Niu, G., Xu, M., Hu, W., Tsang, I., Sugiyama, M., 2018. Co-teaching: Robust training of deep neural networks with extremely noisy labels, in: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems*, Curran Associates, Inc.
- Harrysson, I.J., Cook, J., Sirimanna, P., Feldman, L.S., Darzi, A., Aggarwal, R., 2014. Systematic review of learning curves for minimally invasive abdominal surgery: a review of the methodology of data collection, depiction of outcomes, and statistical analysis. *Annals of surgery* 260, 37–45.
- Hasan, S.K., Linte, C.A., 2019. U-netplus: A modified encoder-decoder u-net architecture for semantic and instance segmentation of surgical instruments from laparoscopic images, in: *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, IEEE. pp. 7205–7211.
- Islam, M., Vibashan, V., Lim, C.M., Ren, H., 2021. St-ntl: Spatio-temporal multitask learning model to predict scanpath while tracking instruments in robotic surgery. *Medical Image Analysis* 67, 101837.
- Jiang, L., Zhou, Z., Leung, T., Li, L.J., Fei-Fei, L., 2018. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels, in: *International Conference on Machine Learning*, PMLR. pp. 2304–2313.
- Jin, Y., Cheng, K., Dou, Q., Heng, P.A., 2019a. Incorporating temporal prior from motion flow for instrument segmentation in minimally invasive surgery video, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 440–448.
- Jin, Y., Cheng, K., Dou, Q., Heng, P.A., 2019b. Incorporating temporal prior from motion flow for instrument segmentation in minimally invasive surgery video github repository. URL: <https://github.com/keyuncheng/MF-TAPNet>.
- Kalia, M., Aleef, T.A., Navab, N., Black, P., Salcudean, S.E., 2021. Co-generation and segmentation for generalized surgical instrument segmentation on unlabelled data, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 403–412.
- Krizhevsky, A., Nair, V., Hinton, G., 2009. Cifar-10 and cifar-100 datasets. URL: <https://www.cs.toronto.edu/kriz/cifar.html> 6, 1.
- Kurmann, T., Márquez-Neila, P., Allan, M., Wolf, S., Sznitman, R., 2021. Mask then classify: multi-instance segmentation for surgical instruments. *International journal of computer assisted radiology and surgery* 16, 1227–1236.
- Laina, I., Rieke, N., Rupperecht, C., Vizcaíno, J.P., Eslami, A., Tombari, F., Navab, N., 2017. Concurrent segmentation and localization for tracking of surgical instruments, in: *International conference on medical image computing and computer-assisted intervention*, Springer. pp. 664–672.
- LeCun, Y., 1998. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/> .
- Lee, K.H., He, X., Zhang, L., Yang, L., 2018. Cleannet: Transfer learning for scalable image classifier training with label noise, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5447–5456.
- Li, F., Kim, T., Humayun, A., Tsai, D., Rehg, J.M., 2013. Video segmentation by tracking many figure-ground segments, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2192–2199.
- Li, W., Wang, L., Li, W., Agustsson, E., Van Gool, L., 2017. Webvision database: Visual learning and understanding from web data. arXiv preprint arXiv:1708.02862 .
- Liu, D., Wei, Y., Jiang, T., Wang, Y., Miao, R., Shan, F., Li, Z., 2020a. Unsupervised surgical instrument segmentation via anchor generation and semantic diffusion, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 657–667.
- Liu, D., Wei, Y., Jiang, T., Wang, Y., Miao, R., Shan, F., Li, Z., 2020b. Unsupervised surgical instrument segmentation via anchor generation and semantic diffusion github repository. URL: <https://github.com/Finspire13/AGSD-Surgical-Instrument-Segmentation>.
- Long, Y., Li, Z., Yee, C.H., Ng, C.F., Taylor, R.H., Unberath, M., Dou, Q., 2021. E-dssr: efficient dynamic surgical scene reconstruction with transformer-based stereoscopic depth perception, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 415–425.
- Ma, X., Huang, H., Wang, Y., Romano, S., Erfani, S., Bailey, J., 2020. Normalized loss functions for deep learning with noisy labels, in: *International Conference on Machine Learning*, PMLR. pp. 6543–6553.

- Mahadevan, S., Athar, A., Ošep, A., Hennen, S., Leal-Taixé, L., Leibe, B., 2020. Making a case for 3d convolutions for object segmentation in videos, in: British Machine Vision Virtual Conference (BMVC).
- Maier-Hein, L., Ross, T., Gröhl, J., Glocker, B., Bodenstedt, S., Stock, C., Heim, E., Götz, M., Wirkert, S., Kenngott, H., et al., 2016. Crowd-algorithm collaboration for large-scale endoscopic image annotation with confidence, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 616–623.
- Marzullo, A., Moccia, S., Catellani, M., Calimeri, F., De Momi, E., 2021. Towards realistic laparoscopic image generation using image-domain translation. *Computer Methods and Programs in Biomedicine* 200, 105834.
- Mascagni, P., Vardazaryan, A., Alapatt, D., Urade, T., Emre, T., Fiorillo, C., Pessaux, P., Mutter, D., Marescaux, J., Costamagna, G., et al., 2021. Artificial intelligence for surgical safety: automatic assessment of the critical view of safety in laparoscopic cholecystectomy using deep learning. *Annals of Surgery*.
- Mayer, N., Ilg, E., Haussler, P., Fischer, P., Cremers, D., Dosovitskiy, A., Brox, T., 2016. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4040–4048.
- Nguyen, T., Mummadi, C., Ngo, T., Beggel, L., Brox, T., 2020. Self: learning to filter noisy labels with self-ensembling, in: International Conference on Learning Representations (ICLR).
- Ni, Z.L., Bian, G.B., Wang, G.A., Zhou, X.H., Hou, Z.G., Xie, X.L., Li, Z., Wang, Y.H., 2020. Barnet: Bilinear attention network with adaptive receptive fields for surgical instrument segmentation, in: Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20, pp. 832–838.
- Nie, D., Gao, Y., Wang, L., Shen, D., 2018. Asdnet: Attention based semi-supervised deep networks for medical image segmentation, in: International conference on medical image computing and computer-assisted intervention, Springer. pp. 370–378.
- Nwoye, C.I., Gonzalez, C., Yu, T., Mascagni, P., Mutter, D., Marescaux, J., Padoy, N., 2020. Recognition of instrument-tissue interactions in endoscopic videos via action triplets, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 364–374.
- Ochs, P., Malik, J., Brox, T., 2013. Segmentation of moving objects by long term video analysis. *IEEE transactions on pattern analysis and machine intelligence* 36, 1187–1200.
- Pakhomov, D., Premachandran, V., Allan, M., Azizian, M., Navab, N., 2019. Deep residual learning for instrument segmentation in robotic surgery, in: International Workshop on Machine Learning in Medical Imaging, Springer. pp. 566–573.
- Pakhomov, D., Shen, W., Navab, N., 2020. Towards unsupervised learning for instrument segmentation in robotic surgery with cycle-consistent adversarial networks, in: 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), IEEE. pp. 8499–8504.
- Perazzi, F., Pont-Tuset, J., McWilliams, B., Van Gool, L., Gross, M., Sorkine-Hornung, A., 2016. A benchmark dataset and evaluation methodology for video object segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 724–732.
- Rieke, N., Tan, D.J., di San Filippo, C.A., Tombari, F., Alsheikhali, M., Belagiannis, V., Eslami, A., Navab, N., 2016. Real-time localization of articulated surgical instruments in retinal microsurgery. *Medical image analysis* 34, 82–100.
- Rolnick, D., Veit, A., Belongie, S., Shavit, N., 2017. Deep learning is robust to massive label noise. *arXiv preprint arXiv:1705.10694*.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical image computing and computer-assisted intervention, Springer. pp. 234–241.
- Ross, T., Zimmerer, D., Vemuri, A., Isensee, F., Wiesenfarth, M., Bodenstedt, S., Both, F., Kessler, P., Wagner, M., Müller, B., et al., 2018. Exploiting the potential of unlabeled endoscopic video data with self-supervised learning. *International journal of computer assisted radiology and surgery* 13, 925–933.
- Sahu, M., Strömsdörfer, R., Mukhopadhyay, A., Zachow, S., 2020. Endo-sim2real: Consistency learning-based domain adaptation for instrument segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 784–794.
- Sawilowsky, S.S., 2009. New effect size rules of thumb. *Journal of modern applied statistical methods* 8, 26.
- Sestini, L., Rosa, B., De Momi, E., Ferrigno, G., Padoy, N., 2021. A kinematic bottleneck approach for pose regression of flexible surgical instruments directly from images. *IEEE Robotics and Automation Letters* 6, 2938–2945.
- Shvets, A.A., Rakhlin, A., Kalinin, A.A., Iglovikov, V.I., 2018a. Automatic instrument segmentation in robot-assisted surgery using deep learning, in: 2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA), IEEE. pp. 624–628.
- Shvets, A.A., Rakhlin, A., Kalinin, A.A., Iglovikov, V.I., 2018b. Automatic instrument segmentation in robot-assisted surgery using deep learning github repository. URL: <https://github.com/ternaus/robot-surgery-segmentation>.
- Song, H., Kim, M., Lee, J.G., 2019. Selfie: Refurbishing unclean samples for robust deep learning, in: International Conference on Machine Learning, PMLR. pp. 5907–5915.
- Song, H., Kim, M., Park, D., Shin, Y., Lee, J.G., 2020. Learning from noisy labels with deep neural networks: A survey. *arXiv preprint arXiv:2007.08199*.
- Sun, D., Yang, X., Liu, M.Y., Kautz, J., 2018. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 8934–8943.
- Takayama, Y., Kaneoka, Y., Maeda, A., Takahashi, T., Uji, M., 2020. Laparoscopic transabdominal preperitoneal repair versus open mesh plug repair for bilateral primary inguinal hernia. *Annals of gastroenterological surgery* 4, 156–162.
- Teed, Z., Deng, J., 2020. Raft: Recurrent all-pairs field transforms for optical flow, in: European conference on computer vision, Springer. pp. 402–419.
- Twinanda, A.P., Shehata, S., Mutter, D., Marescaux, J., De Mathelin, M., Padoy, N., 2016. Endonet: a deep architecture for recognition tasks on laparoscopic videos. *IEEE transactions on medical imaging* 36, 86–97.
- Wang, W., Shen, J., Yang, R., Porikli, F., 2017. Saliency-aware video object segmentation. *IEEE transactions on pattern analysis and machine intelligence* 40, 20–33.
- Wang, Y., Ma, X., Chen, Z., Luo, Y., Yi, J., Bailey, J., 2019. Symmetric cross entropy for robust learning with noisy labels, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 322–330.
- WorldLaparoscopyHospital, . World laparoscopy hospital public repository. <https://www.laparoscopyhospital.com/>.
- Xiao, H., Rasul, K., Vollgraf, R., 2017. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*.
- Xiao, T., Xia, T., Yang, Y., Huang, C., Wang, X., 2015. Learning from massive noisy labeled data for image classification, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2691–2699.
- Yang, C., Lamdouar, H., Lu, E., Zisserman, A., Xie, W., 2021. Self-supervised video object segmentation by motion grouping, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 7177–7188.
- Yang, J., Drake, T., Damianou, A., Maarek, Y., 2018. Leveraging crowdsourcing data for deep active learning an application: Learning intents in alexa, in: Proceedings of the 2018 World Wide Web Conference, pp. 23–32.
- Yang, Y., Loquercio, A., Scaramuzza, D., Soatto, S., 2019a. Unsupervised moving object detection via contextual information separation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 879–888.
- Yang, Y., Loquercio, A., Scaramuzza, D., Soatto, S., 2019b. Unsupervised moving object detection via contextual information separation github repository. URL: [https://github.com/antonilo/unsupervised\\_detection](https://github.com/antonilo/unsupervised_detection).
- Yu, L., Wang, S., Li, X., Fu, C.W., Heng, P.A., 2019. Uncertainty-aware self-ensembling model for semi-supervised 3d left atrium segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 605–613.
- Zhang, Z., Sabuncu, M.R., 2018. Generalized cross entropy loss for training deep neural networks with noisy labels, in: 32nd Conference on Neural Information Processing Systems (NeurIPS).
- Zhao, Z., Jin, Y., Gao, X., Dou, Q., Heng, P.A., 2020. Learning motion flows for semi-supervised instrument segmentation from robotic surgical video, in: International Conference on Medical Image Computing and Computer-

Assisted Intervention, Springer. pp. 679–689.

Zhu, J.Y., Park, T., Isola, P., Efros, A.A., 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks, in: Proceedings of the IEEE international conference on computer vision, pp. 2223–2232.

Zia, A., Essa, I., 2018. Automated surgical skill assessment in rmis training. International journal of computer assisted radiology and surgery 13, 731–739.

Zorn, L., Nageotte, F., Zanne, P., Legner, A., Dallemagne, B., Marescaux, J., de Mathelin, M., 2017. A novel telemanipulated robotic assistant for surgical endoscopy: preclinical application to esd. IEEE Transactions on Biomedical Engineering 65, 797–808.

## Appendix A. Implementation Details

In our implementation, all the segmentation models have a U-Net-like architecture. The *Teacher* network has a 5-convolutional-layers encoder (Figure A1); the *Proxy* network has a 11-convolutional-layers encoder (Figure A2); the *Student* network has the same architecture as TerausNet-16 (Shvets et al. (2018a)), using a VGG-16 architecture as encoder, initialized from ImageNet pre-training. The generator network  $G$  also has a U-Net-like architecture, but uses bilinear-upsampling instead of deconvolution in the expanding path (Figure A3). The discriminator model is implemented using two separate neural networks, one producing a single score as output, another one producing a 16x16 local score-map, in charge of global and local appearance, respectively (Figure A4).

Training parameters, determined from preliminary experiments on external data (*phantom* dataset from Sestini et al. (2021)), are reported in Table A1.

Layer	Input	Kernel	Stride	Filters	Output size	Norm	Activation
Input $E^{OF}(x_t, x_{t+1}) / m^{OF}$	-	-	-	2	256	-	-
Convolution (C1)	$E^{OF}(x_t, x_{t+1}) / m^{OF}$	3	1	32	256	Group (32)	LeakyReLU (0.2)
Max-Pooling (MP1)	C1	2	2	32	128	-	-
Convolution (C2)	MP1	3	1	64	128	Group (32)	LeakyReLU (0.2)
Max-Pooling (MP2)	C2	2	2	64	64	-	-
Convolution (C3)	MP2	3	1	128	64	Group (32)	LeakyReLU (0.2)
Max-Pooling (MP3)	C3	2	2	128	32	-	-
Convolution (C4)	MP3	3	1	256	32	Group (32)	LeakyReLU (0.2)
Max-Pooling (MP4)	C4	2	2	256	16	-	-
Convolution (C5)	MP4	3	1	256	16	Group (32)	LeakyReLU (0.2)
Max-Pooling (MP5)	C5	2	2	256	8	-	-
Convolution (C6)	MP5	3	2	512	4	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.1)	C6	3	1/2	256	8	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.2)	[FSC.1, MP5]	3	1/2	256	16	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.3)	[FSC.2, MP4]	3	1/2	128	32	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.4)	[FSC.3, MP3]	3	1/2	64	64	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.5)	[FSC.4, MP2]	3	1/2	32	128	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.6)	[FSC.5, MP1]	3	1/2	32	256	Group (32)	LeakyReLU (0.2)
Convolution	FSC.6	3	1	1	256	-	Sigmoid

Fig. A1. Network architecture of *Teacher* optical-flow segmentation model.

Layer	Input	Kernel	Stride	Filters	Output size	Norm	Activation
Input $x_t$	-	-	-	3	256	-	-
Convolution (C1)	Input $x_t$	3	2	64	128	Group (32)	LeakyReLU (0.2)
Convolution (C1.2)	C1	3	1	64	128	Group (32)	LeakyReLU (0.2)
Convolution (C2)	C1.2	3	2	128	64	Group (32)	LeakyReLU (0.2)
Convolution (C2.2)	C2	3	1	128	64	Group (32)	LeakyReLU (0.2)
Convolution (C3)	C2.2	3	2	256	32	Group (32)	LeakyReLU (0.2)
Convolution (C3.2)	C3	3	1	256	32	Group (32)	LeakyReLU (0.2)
Convolution (C4)	C3.2	3	2	512	16	Group (32)	LeakyReLU (0.2)
Convolution (C4.2)	C4	3	1	512	16	Group (32)	LeakyReLU (0.2)
Convolution (C5)	C4.2	3	2	512	8	Group (32)	LeakyReLU (0.2)
Convolution (C5.2)	C5	3	1	512	8	Group (32)	LeakyReLU (0.2)
Convolution (C6)	C5.2	3	1	512	8	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.1)	C6	3	1/2	512	16	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.2)	[FSC.1, C4.2]	3	1/2	256	32	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.3)	[FSC.2, C3.2]	3	1/2	128	64	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv. (FSC.4)	[FSC.3, C2.2]	3	1/2	64	128	Group (32)	LeakyReLU (0.2)
Fract. Strided Conv.	[FSC.4, C1.2]	3	1/2	1	256	-	Sigmoid

Fig. A2. Network architecture of *Proxy* segmentation model.

## Appendix B. Additional Qualitative Results

We report additional qualitative results for surgical tool segmentation experiment on EndoVis2017VOS and STRAS datasets (Tables 2&4 in the manuscript), randomly drawn from best and worst 10% predictions of the experiments according to the IoU metric. Results are shown in Figure A5. We also report additional qualitative results for optical-flow

Layer	Input	Kernel	Stride	Filters	Output size	Norm	Activation
Resize Noise ( $n_r$ )	$n [1, 1, 32]$	-	-	32	256	-	-
Input $\text{conc}(x_t, n_r)$	$x_t, n_r$	-	-	3+32	256	-	-
Convolution (C1)	Input $\text{conc}(x_t, n_r)$	3	2	64	128	Layer	LeakyReLU (0.2)
Convolution (C2)	C1	3	2	128	64	Layer	LeakyReLU (0.2)
Convolution (C3)	C2	3	2	256	32	Layer	LeakyReLU (0.2)
Convolution (C4)	C3	3	2	256	16	Layer	LeakyReLU (0.2)
Convolution (C5)	C4	3	2	512	8	Layer	LeakyReLU (0.2)
Convolution (C6)	C5	3	2	512	4	Layer	LeakyReLU (0.2)
Convolution (C,1U)	C6	3	1	512	4	Layer	LeakyReLU (0.2)
Bilinear Upsampling (BU1)	C,1U	-	-	512	8	-	-
Convolution (C,2U)	[BU1, C,1U]	3	1	256	8	Layer	LeakyReLU (0.2)
Bilinear Upsampling (BU2)	C,2U	-	-	256	16	-	-
Convolution (C,3U)	[BU2, C,2U]	3	1	256	16	Layer	LeakyReLU (0.2)
Bilinear Upsampling (BU3)	C,3U	-	-	256	32	-	-
Convolution (C,4U)	[BU3, C,3U]	3	1	128	32	Layer	LeakyReLU (0.2)
Bilinear Upsampling (BU4)	C,4U	-	-	128	64	-	-
Convolution (C,5U)	[BU4, C,4U]	3	1	64	64	Layer	LeakyReLU (0.2)
Bilinear Upsampling (BU5)	C,5U	-	-	64	128	-	-
Convolution (C,6U)	[BU5, C,5U]	3	1	32	128	Layer	LeakyReLU (0.2)
Bilinear Upsampling (BU6)	C,6U	-	-	32	256	-	-
Convolution	BU6	3	1	2	256	-	-

Fig. A3. Network architecture of optical-flow generator model ( $G$ ).

Layer	Input	Kernel	Stride	Filters	Output size	Norm	Activation
Input $E^{OF}(x_t, x_{t+1}) / m^{OF}$	-	-	-	2	256	-	-
Convolution (C1)	$E^{OF}(x_t, x_{t+1}) / m^{OF}$	3	2	64	128	Group (32)	LeakyReLU (0.2)
Convolution (C1.2)	C1	3	1	64	128	-	-
Convolution (C1_shortcut)	$E^{OF}(x_t, x_{t+1}) / m^{OF}$	3	2	64	128	-	-
Residual (R1)	C1_shortcut + C1.2	-	-	64	128	Batch	LeakyReLU (0.2)
Convolution (C2)	R1	3	2	128	64	Group (32)	-
Convolution (C2.2)	C2	3	1	128	64	-	-
Convolution (C2_shortcut)	R1	3	2	128	64	-	-
Residual (R2)	C2_shortcut + C2.2	-	-	128	64	Batch	LeakyReLU (0.2)
Convolution (C3)	R2	3	2	256	32	Batch	LeakyReLU (0.2)
Convolution (C3.2)	C3	3	1	256	32	-	-
Convolution (C3_shortcut)	R2	3	2	256	32	-	-
Residual (R3)	C3_shortcut + C3.2	-	-	256	32	Batch	LeakyReLU (0.2)
Convolution (C4)	R3	3	2	512	16	Batch	LeakyReLU (0.2)
Convolution (C4.2)	C4	3	1	512	16	-	-
Convolution (C4_shortcut)	R3	3	2	512	16	-	-
Residual (R4)	C4_shortcut + C4.2	-	-	512	16	Batch	LeakyReLU (0.2)
Max-Pool (MP)	R4	2	2	512	8	-	-
Flatten+FullyConnected	MP	-	-	1	1	-	Sigmoid

Layer	Input	Kernel	Stride	Filters	Output size	Norm	Activation
Input $E^{OF}(x_t, x_{t+1}) / m^{OF}$	-	-	-	2	256	-	-
Convolution (C1)	$E^{OF}(x_t, x_{t+1}) / m^{OF}$	3	2	32	256	Group (32)	LeakyReLU (0.2)
Max-Pooling (MP1)	C1	3	1	32	128	-	-
Convolution (C2)	MP1	3	1	64	128	Group (32)	LeakyReLU (0.2)
Max-Pooling (MP2)	C2	2	2	64	64	-	-
Convolution (C3)	MP2	3	1	128	64	Group (32)	LeakyReLU (0.2)
Max-Pooling (MP3)	C3	2	2	128	32	-	-
Convolution (C4)	MP3	3	1	256	32	Group (32)	LeakyReLU (0.2)
Max-Pooling (MP4)	C4	2	2	256	16	-	-
Convolution (C5)	MP4	4	1	256	16	Group (32)	LeakyReLU (0.2)
Convolution (C6)	C5	4	1	128	16	Group (32)	LeakyReLU (0.2)
Convolution	C6	4	1	1	16	-	Sigmoid

Fig. A4. Network architecture of discriminator model  $D$ . Top: global discriminator, outputting a single global score; bottom: patch discriminator, outputting a 16x16 score-map.

$n_{\text{epochs}}$	40&40
Batch size	16
LR <sub>GAN</sub>	$3 \times 10^{-3}$
LR <sub>Teacher</sub>	$2 \times 10^{-3}$
LR <sub>Proxy</sub>	$5 \times 10^{-4}$ ( $\div 2 / 5$ epochs, after epoch 20)
LR <sub>Student</sub>	$5 \times 10^{-5}$ ( $\div 2 / 5$ epochs, after epoch 20)
$\beta_1$ <sub>GAN</sub>	0.5
$\beta_2$ <sub>GAN</sub>	0.9
$\beta_1$ <sub>Teacher,Proxy,Student</sub>	0.9
$\beta_2$ <sub>Teacher,Proxy,Student</sub>	0.999
$\epsilon_T = \epsilon_P$	0.5
$\alpha_P = \alpha_S = \alpha$	0.8

Table A1. Training hyper-parameters used in our experiments. Parameters reported: number of training epochs ( $n_{\text{epochs}}$ ) for step-1 (*Teacher* and *Proxy* training) & step-2 (*Student* training), batch size, learning-rates (LR),  $\beta_1$  and  $\beta_2$  for Adam optimizers, *Teacher* and *Proxy* binarization thresholds ( $\epsilon_T$ ,  $\epsilon_P$ ), loss balancing coefficients ( $\alpha_P$ ,  $\alpha_S$ ).

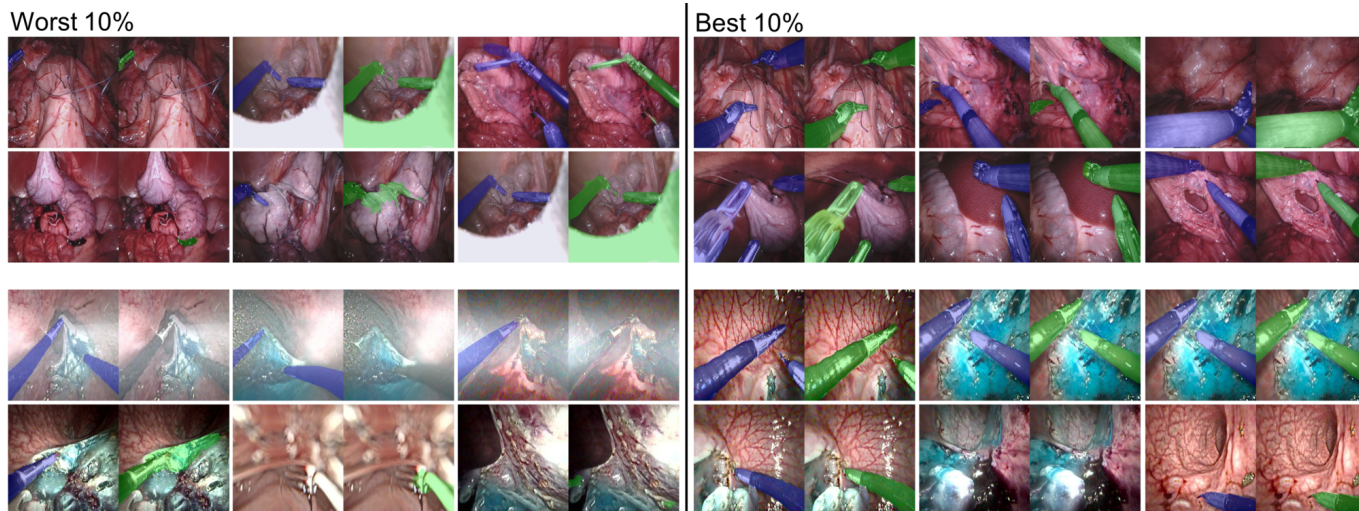


Fig. A5. Qualitative results randomly drawn from worst (left) and best (right) 10% predictions, according to IoU metric on EndoVis2017VOS (top) and STRAS (bottom) datasets. Original frame overlapped with ground-truth (blue) and *Student* network's prediction (green).

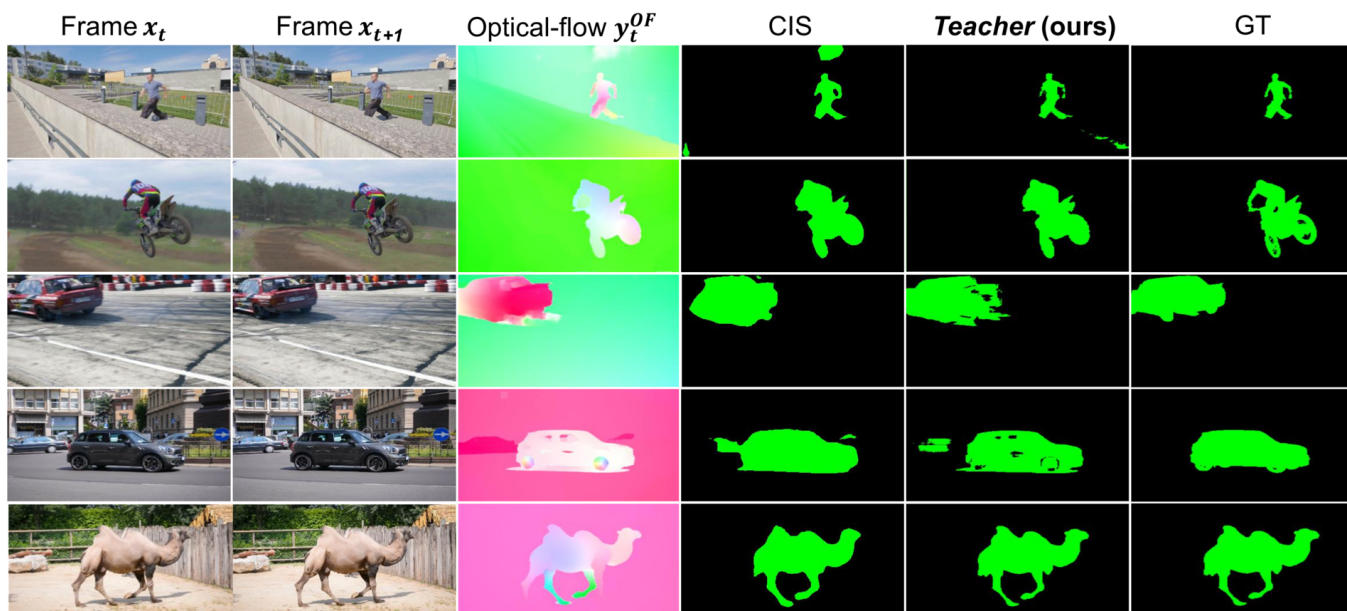


Fig. A6. Optical-flow object segmentation on DAVIS2016 dataset. Qualitative results showing the two frames used for optical-flow computation, optical-flow image after HSV standard conversion, CIS (Yang et al. (2019a)) and *Teacher* (using SegTrackV2 *shape-priors*) predictions, and ground-truth (GT).

object segmentation on DAVIS2016 dataset (Figure A6) for the state-of-the-art CIS approach and our *Teacher* model, trained using SegTrackV2 as *shape-priors*.