# Segmenting the Uterus in Monocular Laparoscopic Images without Manual Input

Toby Collins, Adrien Bartoli, Nicolas Bourdel, and Michel Canis

ALCoV-ISIT, UMR 6284 CNRS/Université d'Auvergne, Clermont-Ferrand, France

**Abstract.** Automatically segmenting organs in monocular laparoscopic images is an important and challenging research objective in computer-assisted intervention. For the uterus this is difficult because of high inter-patient variability in tissue appearance and low-contrast boundaries with the surrounding peritoneum. We present a framework to segment the uterus which is completely automatic, requires only a single monocular image, and does not require a 3D model. Our idea is to use a patient-independent uterus detector to roughly localize the organ, which is then used as a supervisor to train a patient-specific organ segmenter. The segmenter uses a physically-motivated organ boundary model designed specifically for illumination in laparoscopy, which is fast to compute and gives strong segmentation constraints. Our segmenter uses a lightweight CRF that is solved quickly and globally with a single graphcut. On a dataset of 220 images our method obtains a mean DICE score of 92.9%.

## 1 Introduction and Background

The problem of segmenting organs in monocular laparoscopic images without any manual input is important yet unsolved for computer assisted laparoscopic surgery. This is challenging due to multiple factors including inter and intra-patient tissue appearance variability, low-contrast and/or ambiguous organ boundaries, texture inhomogeneity, bleeding, motion blur, partial views, surgical intervention and lens smears. In previous works a manual operator has been needed to identify the organ in one or more training images [3,11]. From these images, models of patient-specific tissue appearance can be learned and used to segment the organ in other images. We present the first methodology to accurately segment an organ in laparosurgery *without any manual input.* Our solution is simple, fast and does not require separate training images, since training and segmentation is performed on the same image. We also do not require patient-specific prior knowledge such as a pre-operative 3D model. Using a 3D model requires registration [11] to give the segmentation (*i.e. segmentation-by-registration*). This shifts the problem burden to registration, which itself is hard to do automatically and reliably for soft organs and monocular laparoscopes [10]. Our approach uses recent work in patient-generic organ detection in laparoscopic images [13]. It was shown that the uterus can be reliably detected in an image *without* patient specific knowledge using a state-of-the-art 2D Deformable Part Model (DPM) detector [8,15] trained on a uterus image database. The problem

of segmentation however was not considered, which is a fundamentally different problem.

For a given image our goal is to compute the binary label matrix $\mathcal{L}(\mathbf{x}) \in \{0, 1\}$ where $\mathcal{L}(\mathbf{x}) = 1$ means pixel $\mathbf{x}$ is on the organ and $\mathcal{L}(\mathbf{x}) = 0$ means it is not. We refer to these as the foreground and background labels respectively. We propose an energy minimisation-based approach to solve $\mathcal{L}$ that incorporates information from the DPM detector to define the energy function. The function is a submodular discrete Conditional Random Field (CRF) that is globally optimised with a *single* graphcut. Much inspiration has come from graphcut-based interactive image segmentation methods [2,14,12] where manual strokes or bounding boxes are used to guide the segmentation. Instead of user interaction, we do this using information from the DPM detector, which in contrast to user interaction information is inherently uncertain. A second major difference is that most graphcut-based methods for optical images use the contrast-sensitive Ising prior from [2], which encourages segmentation boundaries at strong intensity step-edges (*i.e.* points with strong first-order intensity derivatives). However step-edges do not accurately model the appearance of an organ's boundary in laparoscopic images. We show that far better segmentations are obtained using a physically-motivated *trough-sensitive Ising prior*, which is computed from the response of a positive Laplacian of Gaussian ($\text{LoG}^+$) filter (*i.e.* a LoG filter with negative responses truncated to zero). This encourages segmentation boundaries at points with strongly positive *second-order* intensity derivatives.

## 2   Methodology

*Segmentation pipeline.* The main components of our method are illustrated in Fig. 1, which processes an image in five stages. In stage 1 we detect the presence of the organ with the DPM uterus detector from [13]. We take the detector's highest-confidence detection and if it exceeds the detector's threshold we assume the organ is visible and proceed with segmentation. The highest-confidence detection has an associated bounding box $\mathcal{B}$, which gives a rough localisation of the organ. In stage 2 we use $\mathcal{B}$ to train rough appearance models for the organ and background, which are used in the CRF as colour-based segmentation cues. Similarly to GrabCut [14] we use Gaussian Mixture Models (GMMs) with parameters denoted by $\theta_{fg}$ and $\theta_{bg}$ respectively. However unlike GrabCut, we do not iteratively recompute the GMM parameters and the segmentation. This is because with our organ boundary model, the first segmentation is usually very accurate even if the appearance parameters are not. This has the advantage of reduced computation time since we only perform one graphcut.

In stage 3 we use the detection's bounding box to extract a Region Of Interest (ROI) $\mathcal{R}$ around the organ, and all pixels outside $\mathcal{R}$ are labelled background. This reduces computation time because pixels outside $\mathcal{R}$ are not included in the CRF. One cannot naively set $\mathcal{R}$ as the detection's bounding box because there is no guarantee that it will encompass the whole organ, as seen in Fig. 2, bottom row. We normalise $\mathcal{R}$ to have a default width of 200 pixels, which gives sufficiently high

resolution to accurately segment the uterus. The normalisation step is important because it means the CRF energy is independent of the organ's scale. Therefore we do not need to adapt any parameters depending on the organ's physical size, distance to the camera or camera focal length. In stage 4 we construct the CRF which includes information from three important sources. The first is colour information from the foreground and background colour models. The second is edge information from the response of a $\text{LoG}^+$ filter applied to $\mathcal{R}$. The third are spatial priors that give energy to pixels depending on where they are in $\mathcal{R}$. All of the CRF energy terms are submodular which means it can be solved globally and quickly using the maxflow algorithm. In practice this takes between 20-50ms with a standard desktop CPU implementation.
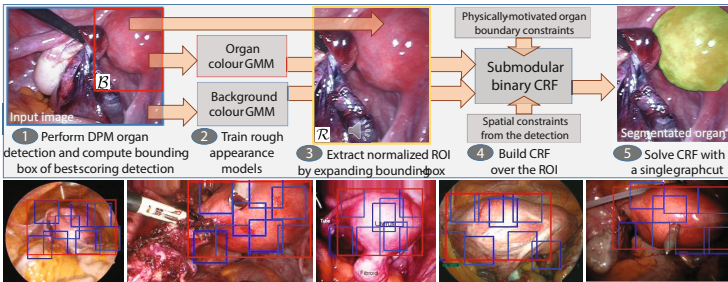


**Fig. 1.** Proposed framework for segmenting the uterus in a monocular laparoscopic image without manual input. The top row shows the five processing stages and the bottom row shows example uterus detections using the DPM detector [13,8].

*The CRF energy function.* The CRF is defined over the ROI $\mathcal{R}$, which is computed by enlarging the bounding box to encompass all likely foreground pixels. This is done by scaling the bounding box about its centre $\mathbf{x}_b$ by a factor of $x\%$. We set this very conservatively to $x = 60\%$, which means all foreground pixels will be within $\mathcal{R}$ when the bounding box of the detection overlaps the ground truth bounding box by at least $\approx 40\%$. In practice we do not normally obtain detections with less than 40% overlap with the ground truth bounding box, because the corresponding detection score would normally be too low to trigger a detection. The CRF energy $E$ is conditioned on $\mathcal{R}$ and $\mathcal{B}$ and is as follows:

$$
\begin{aligned}
E(\mathcal{L}; \mathcal{R}, \mathcal{B}) &\overset{\text{def}}{=} E_{app}(\mathcal{L}; \mathcal{R}) + \lambda_{edge} E_{edge}(\mathcal{L}; \mathcal{R}) + \lambda_{spatial} E_{spatial}(\mathcal{L}; \mathcal{R}, \mathcal{B}) \\
E_{app}(\mathcal{L}; \mathcal{R}) &\overset{\text{def}}{=} \sum_{\mathbf{x} \in \mathcal{R}} \mathcal{L}(\mathbf{x}) E'_{app}(\mathbf{x}; \theta_{fg}) + (1 - \mathcal{L}(\mathbf{x})) E'_{app}(\mathbf{x}; \theta_{bg})
\end{aligned}
\tag{1}
$$

The first term $E_{app}$ denotes the *appearance energy*, which is a standard unary term that encourages pixel labels to agree with the foreground and background GMM models [14]. The term $E'_{app}(\mathbf{x}; \theta)$ denotes the negative density of a GMM parameterised by $\theta$. The terms $E_{edge}$ and $E_{spatial}$ denote the edge and spatial energies, which are unary and pairwise clique energies respectively. The terms $\lambda_{edge}$ and $\lambda_{spatial}$ are weights that govern the relative influence of the energies.
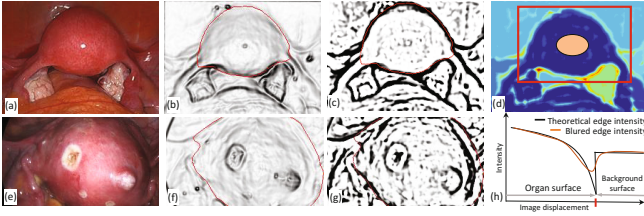
**Fig. 2.** Laparoscopic images of two uteri with different filter response maps (Sobel: (b,f), LoG$^+$: (c,g)), overlaid with manual segmentations. The LoG$^+$ geodesic distance transform $\mathcal{D}$ for (a) is shown in (d), with the detection's bounding box and central ellipse $\mathcal{S}$ overlaid. An illustration of the edge intensity profile across an organ boundary edge is shown in (h).

*A physically-motivated edge energy model based on the LoG$^+$ filter.* The purpose of the edge energy is to encourage a smooth segmentation whose boundary is attracted to probable organ boundaries. In nearly all graphcut-based optical image segmentation methods, this is based on the step-edge model, which says that a transition between labels should occur at regions with high first-order intensity derivatives [2]. However this model does not match well with the physical image formation process in laparoscopic images. This is a combination of the fact that the scene is illuminated by a proximal light source centred close to the camera's optical center, and that because organs are smooth, discontinuities in surface orientation are rare. To see this, consider a point **p** on the organ's boundary with a normal vector **n** in camera coordinates. By definition **n** must be orthogonal to the viewing ray, which implies **n** is approximately orthogonal to the light source vector, so **p** necessarily reflects a very small fraction of direct illumination. Consider now the image intensity profile as we transition from the organ to a background structure (Fig. 2(h)). We observe a smooth intensity fall-off as the boundary is reached, and then a discontinuous jump as we transition to the background. Due to imperfect optics we measure a smooth version of this profile, which is characterised by a smooth intensity trough at a boundary point. *Likely organ boundaries are therefore those image points with strongly positive second-order intensity derivatives*, which can be computed stably with the LoG$^+$ filter. One issue is that edge filters such as LoG$^+$ are also sensitive to superficial texture variation of the organ. An effective way to deal with this is to apply the filter on the red channel only, because red light diffuses deeper into tissue than blue and green light [4]. Fig. 2 illustrates the effectiveness of the LoG$^+$ filter for revealing the uterus boundaries, which we compare to the Sobel step-edge filter.

We define $E_{edge}$ in a similar manner to [2] but replace the intensity difference term by the LoG$^+$ response at the midpoint of two neighbouring pixels **x** and **y**:

$$E_{edge}(\mathcal{L}) \overset{\text{def}}{=} \sum_{(\mathbf{x},\mathbf{y})\in\mathcal{N}} w_{\mathbf{x},\mathbf{y}}(\mathcal{L}) \exp\left(-\text{LoG}^+((\mathbf{x}+\mathbf{y})/2)/2\sigma\right)$$
$$w_{\mathbf{x},\mathbf{y}}(\mathcal{L}) = \begin{cases} 1/d(\mathbf{x},\mathbf{y}) & \text{if } \mathcal{L}(\mathbf{x}) \neq \mathcal{L}(\mathbf{y}) \\ 0 & \text{otherwise} \end{cases} \tag{2}$$

where $\mathcal{N}$ denotes the set of pixel neighbour pairs (we use the standard 8-way connected neighbours from the pixel grid). The term $w_{\mathbf{x},\mathbf{y}} \in \mathbb{R}$ assigns energy when neighbouring pixels have different labels. The function $d$ gives the Euclidean distance between $\mathbf{x}$ and $\mathbf{y}$, which reduces the influence of neighbours that are further away. Inspired by [2] we set $\sigma$ automatically as the standard deviation of the $\mathrm{LoG}^+$ filter across all pixels in $\mathcal{R}$. The $\mathrm{LoG}^+$ filter has a free parameter $\sigma_N$ that pre-smoothes the image to mitigate noise. We have found that results are not highly sensitive to $\sigma_N$, and in all experiments we use $\sigma_N = 3$ pixels with a filter window of 7 pixels.

*Hard labels and spatial energy.* We assign hard labels to pixels in the image that we are virtually certain of either being on the organ or on the background. The job of this is to prevent complete over or under-segmentation in instances when the organ's appearance is very similar to the background. We assign pixels within a small region around the bounding box center $\mathbf{x}_b$ the foreground label, which is valid because the main body of the uterus is always highly convex. Specifically we define a small elliptical region $\mathcal{S}$ by $\mathbf{x} \in \mathcal{S} \Leftrightarrow s^2(\mathbf{x}-\mathbf{x}_b)^\top \mathrm{diag}(1/w, 1/h)(\mathbf{x}-\mathbf{x}_b) \leq 1$, and assign all pixels in $\mathcal{S}$ the foreground label. This is an ellipse with the same aspect ratio as the bounding box, where $w$ and $h$ are the width and height of the bounding box. The scale of $\mathcal{S}$ is given by $s$, which is not a sensitive parameter and in all experiments we use $s = 0.2$. To prevent complete over-segmentation we assign pixels very far from the bounding box the background label. We do this by padding $\mathcal{R}$ by a small amount by replication (we use 20 pixels), and assign the perimeter of the padded image the background label.

The spatial energy encodes the fact that pixels near the detection's center are more likely to be on the organ. We measure distances to the detection's center in terms of geodesics $\mathcal{D}(\mathbf{x}) : \mathcal{R} \to \mathbb{R}^+$ using the $\mathrm{LoG}^+$ filter response as a local metric. This is fast to compute and more informative than the Euclidean distance because it takes into account probable organ boundaries in the image. We compute $\mathcal{D}(\mathbf{x})$ by measuring the distance of $\mathbf{x}$ to $\mathcal{S}$ using the fast marching method. We give a visualisation of $\mathcal{D}$ for the image in Fig. 2 (a) in Fig. 2 (d), with the central ellipse overlaid in red. Dark blue indicates lower distances, and the darkest shade corresponds to a distance of zero. One can see that for most pixels either on the uterus body, or connected to the uterus body by ligaments or the Fallopian tubes, the distance is zero, because for these points there exists a path in the image to $\mathcal{S}$ that does cross an organ boundary. We therefore propose a very simple spatial energy function, which works by increasing the energy of a pixel $\mathbf{x}$ if it is labelled background and has $\mathcal{D}(\mathbf{x}) = 0$. We do this for all pixels within the detection's bounding box, and define the spatial energy as

$$E_{spatial}(\mathcal{L}; \mathcal{D}, \mathcal{B}) \stackrel{\text{def}}{=} \sum_{\mathbf{x} \in R} \begin{cases} 1 & \text{if } \mathcal{L}(\mathbf{x}) = 0 \text{ and } \mathcal{D}(\mathbf{x}) = 0 \text{ and } \mathbf{x} \in \mathcal{B} \\ 0 & \text{otherwise} \end{cases} \qquad (3)$$

The effect of $E_{spatial}$ is to encourage pixels within the bounding box to be labelled foreground if they can reach the detection's center by a path that does not cross points that are likely to be organ boundaries. To improve the computation speed for $E_{spatial}$ we compute $\mathcal{D}$ on a down-sampled version of $\mathcal{R}$ (by a

factor of two). On a standard desktop PC this means $E_{spatial}$ can be computed in approximately 100 to 200ms without significant impact on accuracy.

## 3   Experimental Results

We have evaluated on a new dataset consisting of 235 uterus images of 126 different individuals, which extends the 39-individual database from [13] (Fig. 3). The dataset includes common difficulties caused by pathological shape, surgical change, partial occlusion, strong light fall-off, low-contrast boundaries and over-saturation. The dataset was gathered from patients at our hospital (12 individuals) and demonstration and tuition images from the web (114 patients). 35.0% of the patients had uteri with pathological shape, caused mostly by uterine fibroids. For each image we computed the best-scoring detection from the uterus detector using the accelerated code of [7]. A detection was considered a true positive if the overlap between the detection's bounding box and the manually-computed bounding box exceeded 55% (which is a typical threshold in object detection literature). In total 220 images had true positive detections. In the other 15 images false positives were caused nearly always by strong tool occlusions. We then segmented all images with true positive detections. Because our method is the first to achieve completely automatic organ segmentation in laparoscopic images, there is not a direct baseline method to compare to. We therefore adapted a number of competitive interactive and seed-based segmentation methods, by replacing manual inputs with the output of the uterus detector. These were as follows. (*i*) *GrabCut-I* [14]: we replaced the user-provided bounding box required in GrabCut with the bounding box from the detection, and replaced hard labels from the user with the same hard labels as described above. (*ii*) *Non-iterative GrabCut* (GrabCut-NI): This was the same as GrabCut-I but terminating after one iteration (*i.e.* the appearance models and segmentation were not iteratively refined). (*ii*) *GrowCut* [15]: we used GrowCut with $\mathcal{S}$ as the foreground seed region and the perimeter of $\mathcal{R}$ as the background seed region. (*ii*) *Edge-based Levelset Region growing* (ELR) [9]: we used a well-known levelset region growing method, using $\mathcal{S}$ as the initial seed region. For GrabCut-I, GrabCut-NI, GrowCut and our method, we tested with RGB and illumination-invariant colourspaces. We found negligible differences between the common illumination-invariant colourspaces, so report results with just one (CrCb). The free parameters of the baseline methods were set by hand to maximise their performance on the dataset. The free parameters of our method ($\lambda_{edge}$ and $\lambda_{spatial}$) were set manually with 20 training images, giving $\lambda_{edge} = 90$ and $\lambda_{spatial} = 7$. The training images were no included in the 220 image dataset and were of different patients. We did not use separate training images for the baseline methods, so we could measure their best possible performance on the dataset.

DICE coefficient boxplots (from Matlab's `boxplot`) and summary statistics are given in Fig. 4. We report *p*-values using the two-sample *t*-test with equal variance. The suffixes (RGB) and (CrCb) indicate running a method with RGB and CrCb colourspaces respectively. We also tested whether our method could

**Fig. 3.** Example images from the test dataset and segmentations from our method.



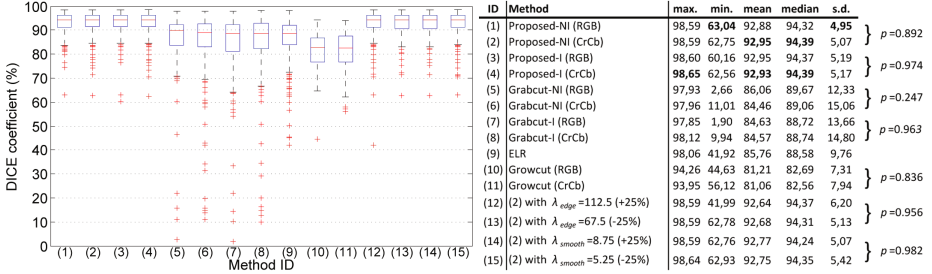| ID | Method | max. | min. | mean | median | s.d. | |
|---|---|---|---|---|---|---|---|
| (1) | Proposed-NI (RGB) | 98,59 | **63,04** | 92,88 | 94,32 | **4,95** | } p =0.892 |
| (2) | Proposed-NI (CrCb) | 98,59 | 62,75 | **92,95** | **94,39** | 5,07 | |
| (3) | Proposed-I (RGB) | 98,60 | 60,16 | 92,95 | 94,37 | 5,19 | } p =0.974 |
| (4) | Proposed-I (CrCb) | **98,65** | 62,56 | 92,93 | **94,39** | 5,17 | |
| (5) | Grabcut-NI (RGB) | 97,93 | 2,66 | 86,06 | 89,67 | 12,33 | } p =0.247 |
| (6) | Grabcut-NI (CrCb) | 97,96 | 11,01 | 84,46 | 89,06 | 15,06 | |
| (7) | Grabcut-I (RGB) | 97,85 | 1,90 | 84,63 | 88,72 | 13,66 | } p =0.963 |
| (8) | Grabcut-I (CrCb) | 98,12 | 9,94 | 84,57 | 88,74 | 14,80 | |
| (9) | ELR | 98,06 | 41,92 | 85,76 | 88,58 | 9,76 | |
| (10) | Growcut (RGB) | 94,26 | 44,63 | 81,21 | 82,69 | 7,31 | } p =0.836 |
| (11) | Growcut (CrCb) | 93,95 | 56,12 | 81,06 | 82,56 | 7,94 | |
| (12) | (2) with $\lambda_{edge}$ =112.5 (+25%) | 98,59 | 41,99 | 92,64 | 94,37 | 6,20 | } p =0.956 |
| (13) | (2) with $\lambda_{edge}$ =67.5 (-25%) | 98,59 | 62,78 | 92,68 | 94,31 | 5,13 | |
| (14) | (2) with $\lambda_{smooth}$ =8.75 (+25%) | 98,59 | 62,76 | 92,77 | 94,24 | 5,07 | } p =0.982 |
| (15) | (2) with $\lambda_{smooth}$ =5.25 (-25%) | 98,64 | 62,93 | 92,75 | 94,35 | 5,42 | |

**Fig. 4.** DICE performance statistics of our proposed method in four configurations (1-4), baseline methods (5-11) and a sensitivity analysis of our method (12-15).

be improved by iteratively retraining the appearance models and resegmenting in the same way as GrabCut (denoted by Proposed-I). Finally, we included a sensitivity analysis of our method, by computing results with $\lambda_{edge}$ and $\lambda_{smooth}$ perturbed from the default by $\pm25\%$. We observe the following. The best performing configurations across all statistics are from the proposed method. There are virtually no differences between our method using RGB or CrCb colourspace, which indicates shading variation does not significantly affect segmentation accuracy. There is also no improvement in our method by iteratively updating the appearance models and resegmenting (Proposed (RGB): $p = 0.998$, Proposed (CrCb): $p = 0.941$). We also see that our method is very stable to a considerable perturbation of the parameters. Fig. 3 shows visually the segmentations from our method (Proposed-NI (CrCb)). The images on the far right show two failure cases. These were caused by a tool occlusion that completely bisected the uterus and a uterus significantly occluded by the laparoscope's optic ring.

## 4   Conclusion

We have presented a method for segmenting the uterus in monocular laparoscopic images that requires no manual input and no patient-specific prior knowledge. We have achieved this using a patient-independent uterus detector to supervise the training of a CRF-based patient-specific segmenter. High accuracy and speed has been obtained by using a physically-motivated organ boundary model based on the $LoG^+$ filter. There are several directions for future work. Firstly, we will

transfer many functions, such as training the GMMs and evaluating the graph constraints onto the GPU for realtime computation. Secondly we will investigate combining our method with a tool segmentation method such as [1]. In terms of applications, our method can be used as a module for automatic laparoscopic video parsing and content retrieval, and for solving problems that have previously required manual organ segmentation. These include building 3D organ models *invivo* [5] and inter-modal organ registration using occluding contours [6].

# References

1. Allan, M., Thompson, S., Clarkson, M.J., Ourselin, S., Hawkes, D.J., Kelly, J., Stoyanov, D.: 2D-3D pose tracking of rigid instruments in minimally invasive surgery. In: Stoyanov, D., Collins, D.L., Sakuma, I., Abolmaesumi, P., Jannin, P. (eds.) IPCAI 2014. LNCS, vol. 8498, pp. 1–10. Springer, Heidelberg (2014)
2. Boykov, Y., Jolly, M.-P.: Interactive graph cuts for optimal boundary amp; region segmentation of objects in N-D images. In: ICCV (2001)
3. Chhatkuli, A., Malti, A., Bartoli, A., Collins, T.: Monocular live image parsing in uterine laparoscopy. In: ISBI (2014)
4. Collins, T., Bartoli, A.: Towards live monocular 3D laparoscopy using shading and specularity information. In: Abolmaesumi, P., Joskowicz, L., Navab, N., Jannin, P. (eds.) IPCAI 2012. LNCS, vol. 7330, pp. 11–21. Springer, Heidelberg (2012)
5. Collins, T., Pizarro, D., Bartoli, A., Canis, M., Bourdel, N.: Realtime wide-baseline registration of the uterus in laparoscopic videos using multiple texture maps. In: MIAR (2013)
6. Collins, T., Pizarro, D., Bartoli, A., Canis, M., Bourdel, N.: Computer-assisted laparoscopic myomectomy by augmenting the uterus with pre-operative MRI data. In: ISMAR (2014)
7. Dubout, C., Fleuret, F.: Exact acceleration of linear object detectors. In: Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., Schmid, C. (eds.) ECCV 2012, Part III. LNCS, vol. 7574, pp. 301–311. Springer, Heidelberg (2012)
8. Felzenszwalb, P., Girshick, R., McAllester, D., Ramanan, D.: Object detection with discriminatively trained part-based models. IEEE PAMI (2010)
9. Li, C., Xu, C., Gui, C., Fox, M.D.: Distance regularized level set evolution and its application to image segmentation. IEEE Trans. Image Process. (2010)
10. Malti, A., Bartoli, A., Collins, T.: Template-based conformal shape-from-motion-and-shading for laparoscopy. In: Abolmaesumi, P., Joskowicz, L., Navab, N., Jannin, P. (eds.) IPCAI 2012. LNCS, vol. 7330, pp. 1–10. Springer, Heidelberg (2012)
11. Nosrati, M., Peyrat, J.-M., Abi-Nahed, J., Al-Alao, O., Al-Ansari, A., Abugharbieh, R., Hamarneh, G.: Efficient multi-organ segmentation in multi-view endoscopic videos using pre-op priors. In: Golland, P., Hata, N., Barillot, C., Hornegger, J., Howe, R. (eds.) MICCAI 2014, Part II. LNCS, vol. 8674, Springer, Heidelberg (2014)
12. Price, B.L., Morse, B.S., Cohen, S.: Geodesic graph cut for interactive image segmentation. In: CVPR (2010)

13. Prokopetc, K., Collins, T., Bartoli, A.: Automatic detection of the uterus and fallopian tube junctions in laparoscopic images. In: Ourselin, S., Alexander, D.C., Westin, C.-F., Cardoso, M.J. (eds.) IPMI 2015. LNCS, vol. 9123, pp. 552–563. Springer, Heidelberg (2015)
14. Rother, C., Kolmogorov, V., Blake, A.: Grabcut - Interactive foreground extraction using iterated graph cuts. ACM Transactions on Graphics (2004)
15. Vezhnevets, V., Konushin, V.: Growcut - Interactive multi-label n-d image segmentation by cellular automata. In: GraphiCon (2005)