

ToF Meets RGB: Novel Multi-Sensor Super-Resolution for Hybrid 3-D Endoscopy

Thomas Köhler^{1,2}, Sven Haase¹, Sebastian Bauer¹, Jakob Wasza¹,
Thomas Kilgus³, Lena Maier-Hein³, Hubertus Feußner⁴,
and Joachim Hornegger^{1,2}

¹ Pattern Recognition Lab, Friedrich-Alexander-Universität Erlangen-Nürnberg

² Erlangen Graduate School in Advanced Optical Technologies (SAOT)
{thomas.koehler,sven.haase}@fau.de

³ Div. Medical and Biological Informatics Junior Group: Computer-Assisted Interventions, German Cancer Research Center (DKFZ) Heidelberg

⁴ Minimally Invasive Therapy and Intervention, Technical University of Munich

Abstract. 3-D endoscopy is an evolving field of research with the intention to improve safety and efficiency of minimally invasive surgeries. *Time-of-Flight* (ToF) imaging allows to acquire range data in real-time and has been engineered into a 3-D endoscope in combination with an RGB sensor (640×480 px) as a hybrid imaging system, recently. However, the ToF sensor suffers from a low spatial resolution (64×48 px) and a poor signal-to-noise ratio. In this paper, we propose a novel multi-frame super-resolution framework to improve range images in a ToF/RGB multi-sensor setup. Our approach exploits high-resolution RGB data to estimate subpixel motion used as a cue for range super-resolution. The underlying non-parametric motion model based on optical flow makes the method applicable to endoscopic scenes with arbitrary endoscope movements. The proposed method was evaluated on synthetic and real images. Our approach improves the peak-signal-to-noise ratio by 1.6 dB and structural similarity by 0.02 compared to single-sensor super-resolution.

1 Introduction

In minimally invasive procedures, reconstructing 3-D surfaces offers opportunities for new applications in addition to conventional 2-D endoscopes. This includes collision detection or augmented reality by registration with preoperative planning data [1]. In terms of hardware, during the past years, three technological directions for 3-D endoscopy emerged. (i) Stereoscopy [2] is a passive technique to acquire surface data. The drawback of stereo vision is the computationally demanding correspondence search and the unreliable results in texture-less regions. (ii) Structured light [3] is established as active acquisition technique. Therefore, the device requires the light source and a sensor for data acquisition placed at a certain distance to observe the situs from different perspectives, which is difficult to accomplish in one single endoscope. (iii) Recently, *Time-of-Flight* (ToF) technology was proposed for 3-D endoscopy to obtain range data in real-time (30 Hz) [4]. In a hybrid 3-D endoscope, the ToF sensor is augmented with an RGB camera to acquire range data fused with complementary color images [5].

Sensor fusion provides the surgeon a comprehensive view of a scene and is beneficial for image analysis [6]. However, today’s ToF sensors suffer from a low spatial resolution and a poor signal-to-noise ratio (SNR) compared to color cameras. Thus, improvement of range data is essential to obtain reliable surface information.

Multi-frame super-resolution methods recover a high-resolution (HR) image from multiple low-resolution (LR) frames with known subpixel displacements [7]. Compared to single image upsampling, such techniques also increase the SNR and preserve edges essential for noisy range data. Recently, super-resolution were applied in 2-D endoscopy [8]. Approaches for color images were also adopted to ToF imaging [9]. An application independent challenge is accurate estimation of subpixel displacements having high impact to super-resolution quality [10]. In literature, several robust methods were proposed [11,12]. Here, super-resolution and motion estimation are formulated as joint optimization which is computationally demanding [12] or restricted to simplified motion models such as rigid motion [11] being an invalid assumption for the considered application.

In this paper, we propose a novel super-resolution framework for range data in a multi-sensor setup. Movements of the endoscope held by the surgeon are used as a cue for super-resolution. Our approach is based on sensor fusion of complementary RGB and range data, which is to the best of our knowledge not considered for multi-frame super-resolution yet. Motion is estimated by computing optical flow on RGB data to obtain accurate displacements for range images. This novelty of our method enables robust motion estimation without computationally demanding joint optimization whereas optical flow avoids restrictions of simplified models essential for realistic laparoscopic scenes. To the best of our knowledge, this is also the first application of super-resolution in 3-D endoscopy.

2 Methods

We address the problem of upsampling K LR range images of resolution $M_1 \times M_2$ denoted as $\mathbf{Y}^{(1)} \dots \mathbf{Y}^{(K)}$ and defined on domain Ω_r . For convenience, we denote $\mathbf{Y}^{(k)}$ as vector $\mathbf{y}^{(k)} \in \mathbb{R}^M$ with $M = M_1 \cdot M_2$ by concatenating all pixels. For each $\mathbf{y}^{(k)}$ there exists a $L_1 \times L_2$ color image $\mathbf{C}^{(k)} \equiv \mathbf{c}^{(k)}$ defined on domain Ω_c captured simultaneously. Each $\mathbf{y}^{(k)}$ and $\mathbf{c}^{(k)}$ is related to a reference frame $\mathbf{y}^{(r)}$ and $\mathbf{c}^{(r)}$ by a geometric transformation modeling 3-D displacements.

Our aim is to determine an HR range image $\mathbf{x} \in \mathbb{R}^N$, $N = s^2 \cdot M$ from K LR frames for the magnification factor $s \in \mathbb{R}$. First, we present sensor data fusion of range and RGB images as key idea in our framework. For super-resolution, an established maximum *a-posteriori* (MAP) estimation scheme is employed [7]. Finally, multi-sensor super-resolution is proposed based on the MAP approach and sensor data fusion for guidance in robust motion estimation.

2.1 Sensor Data Fusion

In hybrid ToF/RGB endoscopy, the incoming light is decomposed by a beam splitter in two components: near-infrared light for the ToF sensor and the residual for the RGB sensor (see Fig. 1). Fusion of both modalities can be tackled by

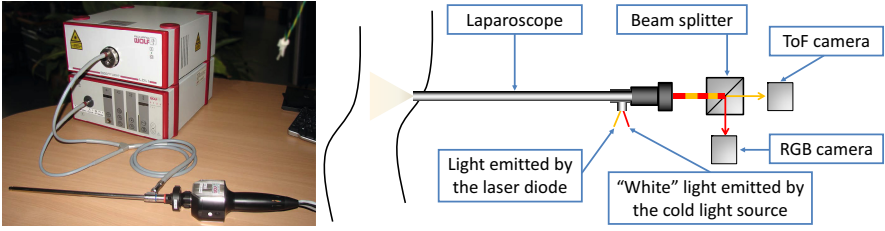


Fig. 1. 3-D endoscope to acquire range and RGB data simultaneously

stereo calibration [5] in general. In our work, we exploit the fact that a beam splitter is used and employ a homographic mapping. For a pair $(\tilde{\mathbf{u}}_r, \tilde{\mathbf{u}}_c)$ of corresponding range and RGB points in homogeneous coordinates, our mapping is given by $\tilde{\mathbf{u}}_c \cong \mathbf{H}_{cr} \tilde{\mathbf{u}}_r$. The homography $\mathbf{H}_{cr} \in \mathbb{R}^{3 \times 3}$ describes pixel-wise alignment and is estimated using a checkerboard calibration pattern with self-encoded markers and least-square estimation as proposed in [5]. Sensor fusion is performed by transforming each $\mathbf{c}^{(k)}$ into the coordinate system of $\mathbf{y}^{(k)}$.

2.2 Maximum A-posteriori Framework

The basic MAP framework [7] is based on a generative image model for mathematical modeling of image acquisition. Super-resolution is implemented by energy minimization based on this model to recover an HR image.

Generative Image Model. The generative image model states the relation between each LR frame $\mathbf{y}^{(k)}$ and the HR image \mathbf{x} to be recovered according to:

$$\mathbf{y}^{(k)} = \gamma_m^{(k)} \mathbf{W}^{(k)} \mathbf{x} + \gamma_a^{(k)} \mathbf{1} + \boldsymbol{\epsilon}^{(k)}. \quad (1)$$

The system matrix $\mathbf{W}^{(k)}$ models geometric displacements between \mathbf{x} and $\mathbf{y}^{(k)}$ as well as blur induced by the camera point spread function (PSF) and downsampling. To take out-of-plane movements and thus diverse range values in successive frames into account, we introduce $\gamma_m^{(k)}$ and $\gamma_a^{(k)}$, where $\mathbf{1} \in \mathbb{R}^M$ denotes the all-one vector. Spatially invariant noise is modeled by $\boldsymbol{\epsilon}^{(k)} \in \mathbb{R}^M$. For a space invariant Gaussian PSF of width σ , the matrix elements are obtained by:

$$W_{mn} = \exp\left(-\|\mathbf{v}_n - \mathbf{u}'_m\|_2^2 / 2\sigma^2\right), \quad (2)$$

where $\mathbf{v}_n \in \mathbb{R}^2$ are the coordinates of the n^{th} pixel in \mathbf{x} and $\mathbf{u}'_m \in \mathbb{R}^2$ are the coordinates of the m^{th} pixel in $\mathbf{y}^{(k)}$ mapped to the HR grid [11]. For efficient memory management, we truncate W_{mn} for $\|\mathbf{v}_n - \mathbf{u}'_m\|_2 > 3\sigma$. Please see our algorithm introduced in section 2.3 for details on the parametrization of this model for range super-resolution proposed in this paper.

MAP Estimator. The objective function to obtain an MAP estimate $\hat{\mathbf{x}}$ for the HR image \mathbf{x} requires a data term and a regularizer weighted by $\lambda > 0$:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \left(\sum_{k=1}^K \left\| \mathbf{y}^{(k)} - \gamma_m^{(k)} \mathbf{W}^{(k)} \mathbf{x} - \gamma_a^{(k)} \mathbf{1} \right\|_2^2 + \lambda \sum_{n=1}^N h_{\tau}((\mathbf{D}\mathbf{x})_n) \right). \quad (3)$$

Algorithm 1. Multi-Sensor Super-Resolution (MSR)

Input: K range images $\mathbf{y}^{(k)}$, RGB data $\mathbf{c}^{(k)}$, reference frame $r = \lceil K/2 \rceil$
Output: Super-resolved range image $\hat{\mathbf{x}}$
for $k = 1 \dots K$ **do**
 $\mathbf{c}^{(k)} := \text{Fuse}(\mathbf{y}^{(k)}, \mathbf{c}^{(k)})$ ▷ see Sect. 2.1
 $\mathbf{w}_c(\mathbf{u}_c) := \text{OpticalFlow}(\mathbf{c}^{(t)}, \mathbf{c}^{(r)})$
 $\mathbf{w}_r(\mathbf{u}_r) := \Delta((l_1 \cdot \mathbf{w}_{c,1}(\mathbf{u}_c) \ l_2 \cdot \mathbf{w}_{c,2}(\mathbf{u}_c))^\top)$ ▷ see Eq. (4)
 $\mathbf{W}^{(k)} := \text{ComposeSystemMatrix}(\mathbf{w}_r(\mathbf{u}_r))$ ▷ see Eq. (2)
 $\gamma_m^{(k)}, \gamma_a^{(k)} := \text{MSAC}(\mathbf{y}^{(r)}, \text{Warp}(\mathbf{y}^{(t)}, \mathbf{w}_r(\mathbf{u}_r)))$
 $\hat{\mathbf{x}}_0 := \text{BicubicUpsampling}(\mathbf{y}^{(r)})$ ▷ initial guess
 $\hat{\mathbf{x}} := \text{SCG}(\hat{\mathbf{x}}_0, \{\mathbf{y}^{(k)}\}, \{\mathbf{W}^{(k)}\}, \{\gamma_m^{(k)}, \gamma_a^{(k)}\})$ ▷ see Eq. (3)

where \mathbf{D} is a high-pass filter and $h_\tau(z) = \tau^2(\sqrt{1 + (z/\tau)^2} - 1)$ is the pseudo Huber loss function used for regularization. For \mathbf{D} we choose a Laplacian to enforce smoothness for \mathbf{x} , which guides the estimation to reliable solutions. However, since the regularizer based on the Huber function penalizes outlier less strictly than a Tikhonov regularization using the L_2 norm, edges are well preserved.

2.3 Multi-Sensor Super-Resolution

In our framework, each $\mathbf{c}^{(k)}$ is aligned to $\mathbf{y}^{(k)}$ after sensor fusion. Motion estimation is performed on RGB images employing optical flow and the displacement fields are projected to the range image domain to compose all system matrices $\mathbf{W}^{(k)}$. The unknown $\gamma_m^{(k)}$ and $\gamma_a^{(k)}$ are determined using robust parameter estimation. We obtain $\hat{\mathbf{x}}$ by minimizing (3) using Scaled Conjugate Gradients (SCG) optimization [13] with a bicubic upsampled version of reference frame $\mathbf{y}^{(r)}$ coincident with $\hat{\mathbf{x}}$ as initial guess. See Algorithm 1 for details of our method.

Optical Flow Estimation. For motion estimation, we determine displacement vector fields $\mathbf{w}_c : \Omega_c \mapsto \mathbb{R}^2$, $\mathbf{w}_c(\mathbf{u}_c) = (\mathbf{w}_{c,1}(\mathbf{u}_c) \ \mathbf{w}_{c,2}(\mathbf{u}_c))^\top$ for RGB images between a reference frame $\mathbf{c}^{(r)}$ and a template $\mathbf{c}^{(t)}$ using optical flow. This transforms each point \mathbf{u}_c from $\mathbf{c}^{(t)}$ to its position \mathbf{u}'_c in $\mathbf{c}^{(r)}$ according to $\mathbf{u}'_c = \mathbf{u}_c + \mathbf{w}_c(\mathbf{u}_c)$. The central frame $\mathbf{c}^{(r)}$ with $r = \lceil K/2 \rceil$ is chosen as reference to minimize the expected displacements between $\mathbf{c}^{(r)}$ and $\mathbf{c}^{(t)}$ for robust flow estimation. Optical flow is computed in a course-to-fine manner using the method proposed by Liu [14]. Once a displacement field \mathbf{w}_c is estimated, it is transformed yielding the range displacement field $\mathbf{w}_r : \Omega_r \mapsto \mathbb{R}^2$:

$$\mathbf{w}_r(\mathbf{u}_r) = \Delta(l_1 \cdot \mathbf{w}_{c,1}(\mathbf{u}_c) \ l_2 \cdot \mathbf{w}_{c,2}(\mathbf{u}_c))^\top, \quad (4)$$

for the resampling operator $\Delta : \mathbb{R}^2 \mapsto \mathbb{R}^2$. We implement Δ as the median of corresponding displacement vectors \mathbf{w}_c in both coordinate directions. To obtain \mathbf{w}_r in the dimension of range data, rescaling by l_i , $0 < l_i \leq 1$ is required, where

l_i denotes the ratio of resolutions between $\mathbf{y}^{(k)}$ and $\mathbf{c}^{(k)}$. Then \mathbf{w}_r is used to compose the system matrices for each frame according to Eq. 2.

Range Diversity Correction. If we allow general 3-D movements of the endoscope such as out-of-plane translation, this results in an offset for range values in successive frames. Neglecting this effect as implicitly done in related super-resolution approaches [9] leads to biased reconstructions. This problem can be mathematically compared to the fusion of intensity images differing photometrically. Therefore, we adopt a photometric registration scheme to range correction. First, two frames to be corrected are assumed to be geometrically aligned by warping them according to the precomputed optical flow displacement field. Let y be a range value in reference frame $\mathbf{y}^{(r)}$. The corresponding range value y' in template frame $\mathbf{y}^{(t)}$ is given according to the affine model $y' = \gamma_m \cdot y + \gamma_a$.

We utilize an M-estimator sample consensus (MSAC) for robust estimation of γ_m and γ_a as suggested by Capel [15] for photometric registration. These parameters are plugged into the generative image model (1) for $\gamma_m^{(k)}$ and $\gamma_a^{(k)}$. For the reference $\mathbf{y}^{(r)}$ we set $\gamma_m^{(r)} = 1$ and $\gamma_a^{(r)} = 0$ to obtain a super-resolved image having the same measurement range as the reference frame.

3 Experiments and Results

We compared multi-sensor super-resolution (MSR) to the conventional single-sensor approach (SSR) where optical flow is estimated on range data. The PSF width was set to $\sigma = 0.5$ and for regularization using Huber function we set $\lambda = 70$ and $\tau = 5 \cdot 10^{-3}$ determined empirically using a grid search. SCG was used with termination tolerance 10^{-3} for pixels of \mathbf{x} and the objective function value. The maximum iteration number was set to 50. Super-resolution was applied with magnification $s = 4$ in a sliding window scheme over time using successive $K = 31$ frames (30 template and one reference frame) per window. This improves the robustness and our method is able to recover from failures caused e. g. by misregistration of single frames in highly dynamic scenes. Supplementary material for our experiments is available on our web page.¹

Synthetic Data. For quantitative assessment, six synthetic data sets based on ground truth data were generated. We used a ToF/RGB simulator to obtain RGB and range data from a model of a laparoscopic scene designed in collaboration with a medical expert (see Fig 2). The resolutions for RGB (640×480 px) and range images (64×48 px) are equal to those of the hybrid 3-D endoscope used in experiments for real data. Each LR frame is a downsampled version of the ground truth and disturbed by a Gaussian PSF ($\sigma_b = 0.5$) as well as additive, zero-mean, Gaussian noise ($\sigma_n = 0.05$). Random motion of the camera was used to simulate movements of the endoscope held by a surgeon. Small displacements of endoscopic tools and organs simulated minimally invasive surgery. As quality metrics we employed the peak-signal-to-noise ratio (PSNR) and structural

¹ <http://www5.cs.fau.de/research/data/>

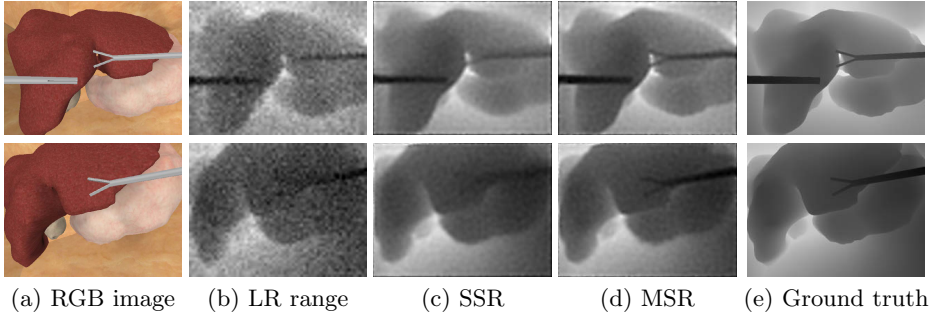


Fig. 2. Synthetic sequences S4 and S5: RGB data (a), LR range data (b), results for SSR (c) and the proposed MSR (d) compared to ground truth data (e)

Table 1. PSNR in dB (SSIM in brackets) for synthetic data. Each result is averaged over 10 sub-sequences per set. We compared bicubic upsampling (second column) to SSR (third column) and our MSR approach with the proposed range correction scheme (fourth column) and without range correction (last column).

Sequence	Interpolation (bicubic)	Range corr.		No corr. MSR
		SSR	MSR	
S1	24.28 (0.57)	28.16 (0.87)	29.85 (0.89)	29.83 (0.89)
S2	25.23 (0.59)	30.78 (0.91)	31.13 (0.92)	30.91 (0.91)
S3	25.83 (0.60)	31.89 (0.92)	32.72 (0.93)	31.93 (0.93)
S4	25.58 (0.59)	29.06 (0.89)	31.17 (0.91)	29.19 (0.91)
S5	26.58 (0.59)	30.36 (0.91)	32.77 (0.93)	31.35 (0.93)
S6	26.26 (0.57)	28.43 (0.89)	30.66 (0.92)	30.28 (0.92)
Mean	25.63 (0.58)	29.78 (0.90)	31.38 (0.92)	30.58 (0.91)

similarity (SSIM). For comparison, we also evaluated bicubic interpolation as a fast and simple upsampling technique. MSR was evaluated with and without range correction to justify our correction scheme. See Tab. 1 for PSNR and SSIM measures averaged over ten subsequent sequences in sliding window processing.

Real Data. For qualitative evaluation, we acquired real data using a hybrid 3-D endoscope prototype manufactured by Richard Wolf GmbH, Knittlingen, Germany. Therefore, a liver phantom and two surgical tools were measured with a frame rate of 30 fps. Range and RGB images were captured and the endoscope was slightly moved during acquisition. Raw data compared to super-resolved data is shown in Fig. 3. See Fig. 4 for a 3-D mesh created for one sequence.

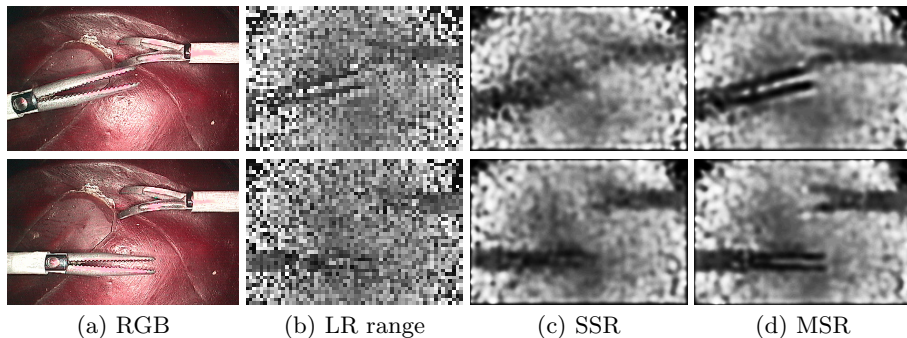


Fig. 3. Liver phantom acquired by a hybrid 3-D endoscope: RGB images (a), LR range data (b) and the results of SSR (c) as well as the proposed MSR (d)

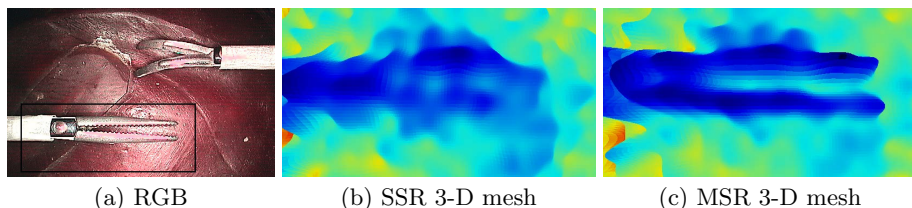


Fig. 4. RGB image (a) and 3-D meshes for SSR (b) and the proposed MSR (c)

4 Discussion

For synthetic images, super-resolution yields more reliable 3-D surfaces compared to bicubic interpolation, especially due to denoising implicitly performed in MAP estimation (see Tab. 1). The proposed MSR approach clearly outperforms SSR indicated by increased mean PSNR (SSIM) of 1.6 dB (0.02). We could also verify improved reconstruction of depth discontinuities by visual inspection, e. g. for surgical tools (see Fig. 2). For sequences containing large out-of-plane movements (S4 and S5), we observed a particularly biased reconstruction if no range correction is applied indicated by decreased PSNR. For real data, the proposed MSR approach recovers the liver surface and endoscopic tools barely visible in raw data as well as in the result of SSR (see Fig. 3).

In our experiments, super-resolution was performed *off-line* to obtain one single HR image from multiple LR frames. Please note, that this limitation may be avoided by dynamic estimation schemes [10] to enable an *on-line* implementation, which is however beyond the scope of this paper.

5 Conclusions

In this paper, a multi-sensor framework for super-resolution in hybrid 3-D endoscopy has been introduced. Range super-resolution is guided by RGB images

acquired simultaneously. Our method is a general-purpose technique to overcome low spatial resolutions and poor SNR of today's ToF sensors. For applications such as segmentation or classification, super-resolution holds great potential and may help to make the breakthrough of ToF imaging in minimally invasive surgery. Beyond ToF/RGB endoscopy, our method may also be applicable in related hybrid range imaging systems, which is part of our future work.

Acknowledgments. The authors gratefully acknowledge funding of the Erlangen Graduate School in Advanced Optical Technologies (SAOT) by the German National Science Foundation (DFG) in the framework of the excellence initiative and the support by the DFG under Grant No. HO 1791/7-1. This research was funded by the Graduate School of Information Science in Health (GSISH) and the TUM Graduate School. We thank the Metrilus GmbH for their support.

References

1. Röhl, S., Bodenstedt, S., Suwelack, S., Kenngott, S., Mueller-Stich, P., Dillmann, R., Speidel, S.: Real-time surface reconstruction from stereo endoscopic images for intraoperative registration. In: Proc. SPIE, vol. 7964, p. 796414 (2011)
2. Field, M., Clarke, D., Strup, S., Seales, W.: Stereo endoscopy as a 3-d measurement tool. In: EMBC 2009, pp. 5748–5751 (2009)
3. Schmalz, C., Forster, F., Schick, A., Angelopoulou, E.: An endoscopic 3d scanner based on structured light. *Med. Image Anal.* 16(5), 1063–1072 (2012)
4. Penne, J., Höller, K., Stürmer, M., Schrauder, T., Schneider, A., Engelbrecht, R., Feußner, H., Schmauss, B., Hornegger, J.: Time-of-Flight 3-D Endoscopy. In: Yang, G.-Z., Hawkes, D., Rueckert, D., Noble, A., Taylor, C. (eds.) MICCAI 2009, Part I. LNCS, vol. 5761, pp. 467–474. Springer, Heidelberg (2009)
5. Haase, S., Forman, C., Kilgus, T., Bammer, R., Maier-Hein, L., Hornegger, J.: ToF/rgb sensor fusion for augmented 3d endoscopy using a fully automatic calibration scheme. In: Tolxdorff, T., Deserno, T.M., Handels, H., Meinzer, H.P. (eds.) *Bildverarbeitung für die Medizin 2012*, pp. 467–474. Springer, Heidelberg (2012)
6. Haase, S., Wasza, J., Kilgus, T., Hornegger, J.: Laparoscopic instrument localization using a 3-d time-of-flight/rgb endoscope. In: WACV 2013, pp. 449–454 (2013)
7. Park, S.C., Park, M.K., Kang, M.G.: Super-resolution image reconstruction: a technical overview. *IEEE Signal Process Mag.* 20(3), 21–36 (2003)
8. De Smet, V., Namboodiri, V., Van Gool, L.: Super-resolution techniques for minimally invasive surgery. In: *Proceedings AE-CAI 2011*, pp. 41–50 (2011)
9. Schuon, S., Theobalt, C., Davis, J., Thrun, S.: Lidarboost: Depth superresolution for tof 3d shape scanning. In: CVPR 2009, pp. 343–350 (2009)
10. Farsiu, S., Robinson, D., Elad, M., Milanfar, P.: Advances and challenges in super-resolution. *Int. J. Imaging Syst. Technol.* 14(2), 47–57 (2004)
11. Tipping, M.E., Bishop, C.M.: Bayesian image super-resolution. In: *Adv. Neural Inf. Process. Syst.*, pp. 1303–1310. MIT Press (2003)
12. Fransens, R., Strecha, C., Van Gool, L.: Optical flow based super-resolution: A probabilistic approach. *Comput. Vis. Image Underst.* 106(1), 106–115 (2007)
13. Nabney, I.T.: NETLAB: Algorithms for Pattern Recognition, 1st edn. *Advances in Pattern Recognition*. Springer (2002)
14. Liu, C.: Beyond Pixels: Exploring New Representations and Applications for Motion Analysis. PhD thesis, Massachusetts Institute of Technology (2009)
15. Capel, D.: Image mosaicing and super-resolution. PhD thesis, Oxford (2004)