# Self-Supervised Generative Adversarial Network for Depth Estimation in Laparoscopic Images

Baoru Huang[1,2], Jianqing Zheng[3], Anh Nguyen[1], David Tuch[4], Kunal Vyas[4], Stamatia Giannarou[1,2], and Daniel S. Elson[1,2]

[1] The Hamlyn Centre for Robotic Surgery, Imperial College London, SW7 2AZ, UK
Baoru.Huang18@imperial.ac.uk
[2] Department of Surgery & Cancer, Imperial College London, SW7 2AZ, UK
[3] The Kennedy Institute of Rheumatology, University of Oxford, UK
[4] Lightpoint Medical Ltd.

**Abstract.** Dense depth estimation and 3D reconstruction of a surgical scene are crucial steps in computer assisted surgery. Recent work has shown that depth estimation from a stereo images pair could be solved with convolutional neural networks. However, most recent depth estimation models were trained on datasets with per-pixel ground truth. Such data is especially rare for laparoscopic imaging, making it hard to apply supervised depth estimation to real surgical applications. To overcome this limitation, we propose SADepth, a new self-supervised depth estimation method based on Generative Adversarial Networks. It consists of an encoder-decoder generator and a discriminator to incorporate geometry constraints during training. Multi-scale outputs from the generator help to solve the local minima caused by the photometric reprojection loss, while the adversarial learning improves the framework generation quality. Extensive experiments on two public datasets show that SADepth outperforms recent state-of-the-art unsupervised methods by a large margin, and reduces the gap between supervised and unsupervised depth estimation in laparoscopic images.

**Keywords:** Depth Estimation · Laparoscopic Images · Generative Adversarial Network

## 1 Introduction

Robot-assisted minimally invasive surgery with stereo laparoscopic vision has become popular due to the advantages of enhanced movement range, precision, vision and proficiency [23,24,33]. Surgical scene depth estimation is a fundamental problem in image-guided intervention and has received substantial prior interest to its promise for robot navigation, 3D registration between pre- and intra-operative organ models, and augmented reality [31]. Obtaining depth maps is not trivial due to the inherent problems such as tissue deformation, specular reflections, and lack of photometric constancy across frames [21].

Several traditional methods used multi-view stereo algorithms such as Simultaneous Localization and Mapping (SLAM) [12] and Structure from Motion

(SfM) [20], but these struggle with less textured tissues. More recently deep learning-based depth estimation has used RGB images as the training data and Convolutional Neural Networks (CNNs) for supervised learning [6,4]. To produce accurate results in less than a second of GPU time, Luo *et al.* [22] treated the problem as a multi-class classification indicating all possible disparities, and exploited a product layer to simplify the representations of a Siamese architecture. Chang *et al.* [2] proposed PSMNet, where the capacity of global context information at different scales and locations could be extracted by a spatial pyramid pooling module to form a cost volume. Duggal *et al.* [5] sped up the runtime of stereo matching and developed a differentiable PatchMatch module that could discard most disparities without the need of full cost volume evaluation.

The methods above are fully supervised and require ground truth depth during training. However, acquiring per-pixel ground truth depth data is challenging for real-world settings [19] and especially for laparosocpic vision where port space is limited, working distance is short and sterilization is required [15]. One alternative is self-supervised training of depth estimation models using image reconstruction as the supervisory signal [7]. The input is usually a set of images in the form of monocular or stereo images [34]. Godard *et al.* [9] proposed a training loss that included a left-right depth consistency term and a reconstruction term for single image depth estimation, despite the absence of ground truth depth. This was extended by [10] with full-resolution multi-scale sampling to reduce visual artifacts, and a minimum reprojection loss to robustly handle occlusions. Johnston *et al.* [18] further closed the gap with fully-supervised methods by including a self-attention mechanism and made use of contextual information. Ye *et al.* [31] proposed a deep learning framework for surgical scene depth estimation in self-supervised mode for scalable data acquisition by adopting a differentiable spatial transformer and an autoencoder.

In this paper, we present a new method for self-supervised adversarial depth estimation: SADepth. A U-Net architecture [27] was adopted as a generative structure and fed with stereo pairs as inputs to benefit from complementary information. To cope with local minima caused by classic photometric reprojection loss, we applied the disparity smoothness loss and formed the network across multiple scales. The use of a generative adversarial network (GAN) allowed us to improve the reconstructed image quality, which formed a supervisory signal for training, while keeping the overall end-to-end optimization objective.

## 2   Methodology

### 2.1   Overview

Here we describe the proposed self-supervised adversarial depth estimation framework, SADepth. Stereo depth estimation predicts depth maps $\boldsymbol{D}^{\mathrm{l}}, \boldsymbol{D}^{\mathrm{r}} \in \mathbb{R}_{+}^{h \times w}$ based on the stereo RGB images $\boldsymbol{I}^{\mathrm{l}}, \boldsymbol{I}^{\mathrm{r}} \in \mathbb{R}_{+}^{h \times w \times 3}$ of height and width $h, w$. A generative network $\mathcal{G}$ with stereo image pairs $\boldsymbol{I}^{\mathrm{l}}$ and $\boldsymbol{I}^{\mathrm{r}}$ as inputs, was used to produce two distinct left and right disparity maps $\boldsymbol{d}^{\mathrm{l}}$ and $\boldsymbol{d}^{\mathrm{r}}$, *i.e.* $\boldsymbol{d}^{\mathrm{l}}, \boldsymbol{d}^{\mathrm{r}} =$

**Fig. 1.** Overview of the self-supervised adversarial depth estimation network, SADepth.

$\mathcal{G}(\boldsymbol{I}^{l}, \boldsymbol{I}^{r})$. As the two disparity maps were generated from different input images, a 'reprojection sampler' [17] could be used for photometric reprojection loss computation of mutual counter-parts, *i.e.* reconstructed left and right images $\boldsymbol{I}^{l*}$ and $\boldsymbol{I}^{r*}$. The discriminator $\mathcal{D}$ was exploited to indicate if the reconstructed images were real or fake (original input images were regarded as real). By forcing the reconstructed image to be consistent with the original input, we could derive accurate disparity maps for depth inference, as shown in the following sections.

## 2.2 Network Architecture

*Generator.* The generator followed the general U-Net [27] architecture consisting of an encoder-decoder network, where the encoder was designed to obtain compact image representations and the decoder produced disparity maps for left and right input images, recovering them at the original scale (illustrated in Figure 3). Encoder-decoder skip connections were applied to represent deep abstract features while preserving local information. To make the model compact - and different from less streamlined previous approaches which had two branches or two sub-networks for the encoder [2] [26] [16] - we first concatenated the left and right images into a 6-channel tensor and then fed it to a ResNet18 model [13]. The input size was $\# \; channels \times h \times w = 6 \times 192 \times 384$. Similar to [9], our decoder was formed of five cascaded blocks where each block had four parts: the first convolutional layer, an upsampling layer, a concatenation manipulation, and the second convolutional layer. In the upsampling layer, features were interpolated to twice the input size and both convolutional layers were followed by an *ELU* activation function [3]. In particular, sigmoids were applied at the output to generate a 2-channel tensor representing the left and right disparity $\mathbf{d}^{l}$ and $\mathbf{d}^{r}$. Finally the sigmoid outputs were converted to depth by $\mathbf{D}^{l(r)} = 1/(a\mathbf{d}^{l(r)} + b)$, where parameters $a$ and $b$ were selected to constrain the depth $\mathbf{D}^{l(r)}$ between 0.1

**Fig. 2.** The detailed architecture of the SADepth generator and discriminator. The generator was an autoencoder architecture with concatenated stereo image pairs as inputs and left and right disparity maps as outputs using a sigmoid function. These outputs were then transformed to reconstruct the counter-part camera input images using a 'reprojection sampler', and these reconstructed images were fed into the discriminator together with the original input image pair. The discriminator output a scalar indicating whether the reconstructed images generated from the 'reprojection sampler' were real or fake.

and 100 units. The depth maps were then back-projected into point clouds by applying the intrinsic parameters and using the counter-part camera's extrinsic parameters to form reconstructed stereo images. The structural similarity between the original and reconstructed images was regarded as a supervisory signal to train the generator (see section 2.3 for the generator loss).

*Discriminator.* Godfellow *et al.* [11] introduced a generative adversarial learning strategy and presented impressive results for image generation tasks. GANs have been widely exploited in different tasks with different GAN models including *e.g.* DualGAN [32] and CycleGAN [35]. To improve the generation quality of the reconstructed images $I^{l*}$ and $I^{r*}$, and following the work in [26] for natural scenes, we applied an adversarial learning strategy for laparoscopic images to include geometry constraints during training and force the network to make a consistent depth map prediction. The original input stereo image pairs and reconstructed images $I^{r*}$ and $I^{l*}$ generated from the 'reprojection sampler' were fed into the discriminator $\mathcal{D}$, which consisted of convolutional, batch normalization and activation function layers and classified the input and reconstructed images as real or fake. As training progressed, the reconstructed images became more similar to the original inputs, while the discriminator also became better at distinguishing between the input and reconstructed images, resulting in an overall improvement of the associated disparity maps.

### 2.3   Training Losses

*Generator Loss.* In the depth estimation generator network $\mathcal{G}$, the loss $\mathcal{L}_{\text{rec}}^{\text{r}}$ was formed from the appearance matching loss $\mathcal{L}_{\text{ap}}^{\text{r}}$ and disparity smoothness loss $\mathcal{L}_{\text{ds}}^{\text{r}}$

$$\mathcal{L}_{\text{rec}}^{\text{r}} = \mathcal{L}_{\text{ap}}^{\text{r}} + \alpha_{\text{ds}}\mathcal{L}_{\text{ds}}^{\text{r}} \tag{1}$$

where $\alpha_{\text{ds}}$ balanced the loss magnitude of the two parts to stabilize the training and was set to 0.001.

*Appearance-Matching Loss.* Self-supervised training typically assumes that the appearance and material properties (*e.g.* brightness and Lambertian) of object surfaces are consistent between frames. A local structure-based appearance loss [9] can effectively improve the depth estimation performance compared with simple pairwise pixel differences [34]. Following [10], we exploited the appearance-matching loss as part of the generator loss which forced the reconstructed image to be similar to the corresponding training inputs. During the training, the right disparity map $\mathbf{d}^{\text{r}}$ generated by the autoencoder was then transformed to produce $\boldsymbol{I}^{\text{r}*}$ – a reconstruction of the original right input image – using RGB intensity information from the counter-part camera image $\boldsymbol{I}^{\text{l}}$ (see Fig. 1). This was achieved by first converting the disparity map $\mathbf{d}^{\text{r}}$ to a depth map $\mathbf{D}^{\text{r}}$, from which a point cloud of the surgical scene could be generated. Then the point cloud was transferred into the other camera's coordinate system and projected onto its image plane. The reconstructed input image $\boldsymbol{I}^{\text{r}*}$ was generated with bilinear interpolation for each output pixel using the weighted sum of the four neighboring intensities. In contrast to [7], this bilinear sampling was locally fully differentiable, which allowed it to be integrated into the fully convolutional architecture without requiring simplification or approximation of the cost function. To compare the reconstructed image $\boldsymbol{I}^{\text{r}*}$ and the original input image $\mathbf{I}^{\text{r}}$, a combination of structural similarity (SSIM) index [28] and $\mathcal{L}_1$ loss were applied as the photometric image reconstruction cost $\mathcal{L}_{ap}^{\text{r}}$:

$$\mathcal{L}_{ap}^{\text{r}} = \frac{1}{N}\sum_{i,j}\frac{\gamma}{2}(1 - \text{SSIM}(I_{ij}^{\text{r}}, I_{ij}^{r*})) + (1-\gamma)\|I_{ij}^{\text{r}} - I_{ij}^{r*}\|_1 \tag{2}$$

where $N$ denotes the number of pixels and $\gamma$ represents the weighting for L1-norm loss term, which was set to 0.85. Similar to [9], the calculation of SSIM here was simplified to a $3 \times 3$ block filter instead of a Gaussian. The training of the depth estimation generator then involved minimizing the reconstruction loss between input and reconstructed images.

*Disparity Smoothness Loss.* Since disparities should be locally smooth and discontinuities usually occur at image gradients, we applied the disparity smoothness loss to penalize unexpected discontinuities in the disparity maps. Following [14], this cost was an edge-aware term weighted with the input image gradients $\partial\mathbf{I}$:

$$\mathcal{L}_{\text{ds}}^{\text{r}} = \frac{1}{N} \sum_{ij} |\partial_x(\mathbf{d}_{ij}^{\text{r}})| e^{-|\partial_x \boldsymbol{I}_{ij}^{\text{r}}|} + |\partial_y(\mathbf{d}_{ij}^{\text{r}})| e^{-|\partial_y \boldsymbol{I}_{ij}^{\text{r}}|} \tag{3}$$

where $\mathbf{d}^{\text{r}}$ represents the generated disparity map and $\boldsymbol{I}^{\text{r}}$ is the original input right image.

*Discriminator Loss.* The adversarial objective of the generative network can be expressed as follows:

$$\mathcal{L}_{\text{gan}}^{\text{r}}(\boldsymbol{I}^{\text{r}}, \boldsymbol{I}^{\text{r}*}; \mathcal{G}, \mathcal{D}) = \mathbb{E}_{\boldsymbol{I}^{\text{r}} \sim P(\boldsymbol{I}^{\text{r}})}[\log(\mathcal{D}(\boldsymbol{I}^{\text{r}}))] + \mathbb{E}_{\boldsymbol{I}^{\text{r}*} \sim P(\boldsymbol{I}^{\text{r}*})}[\log(1 - \mathcal{D}(\boldsymbol{I}^{\text{r}*}))] \tag{4}$$
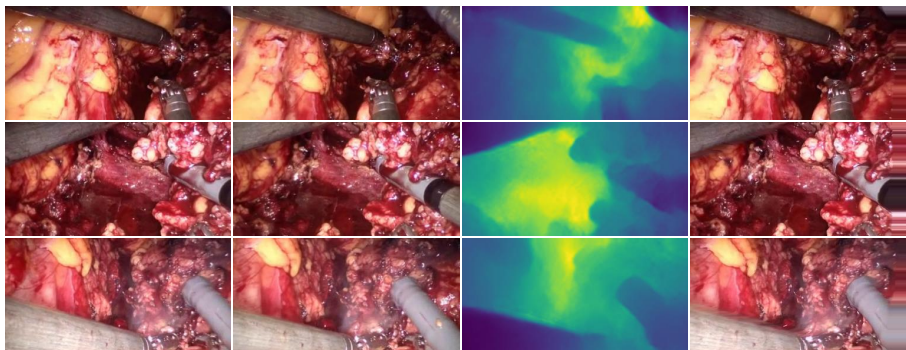
where a cross-entropy loss measured the expectation of the reconstructed image $\boldsymbol{I}^{\text{r}*}$ against the distribution of the input image $\boldsymbol{I}^{\text{r}}$. Note that both generator and discriminator losses included losses for left and right images but only the right image equations are shown.

*Multi-Scale Loss.* One remaining issue with the above learning pipeline was that the training objective risked becoming stuck in local minima due to the application of a photometric reprojection loss [29]. The strategy introduced in [34] indicated that combining the individual losses across multiple scales in the decoder was effective, which could improve the depth estimation performance and reduce sensitivity to architectural choices. Hence, the lower resolution depth maps (from the intermediate layers) were first upsampled to the input image resolution and then reprojected and resampled, with the errors computed at the higher input resolution. This manipulation is similar to matching patches, which enables low-resolution disparity maps to warp an entire patch of pixels in a high resolution image while promoting the depth maps at every scale to reconstruct the high resolution input image as accurately as possible [10].

*Joint Optimization Loss* Finally, the joint optimization loss was a combination of generator loss and adversarial loss, written as:

$$\mathcal{L}_{\text{total}} = \frac{1}{m} \sum_{s=1}^{m} \frac{\mathcal{L}_s^{\text{l}} + \mathcal{L}_s^{\text{r}}}{2} = \frac{1}{m} \sum_{s=1}^{m} \left( \alpha(\mathcal{L}_{\text{rec}}^{\text{l}} + \mathcal{L}_{\text{rec}}^{\text{r}}) + \beta(\mathcal{L}_{\text{gan}}^{\text{l}} + \mathcal{L}_{\text{gan}}^{\text{r}}) \right) \tag{5}$$

**Training** The depth estimation procedure was trained based on the reconstruction supervision signal and no per-pixel depth ground truth labels were needed. The augmentation of input data was performed on the fly by flipping 50 % of the input images horizontally and reorienting the stereo pairs. Parameter $m$ was set to 4, which means that there were 4 output scales with resolutions $\frac{1}{2^0}$, $\frac{1}{2^1}$, $\frac{1}{2^2}$ and $\frac{1}{2^3}$ of the input resolution. $\alpha$ and $\beta$ were set to 0.5.

**Fig. 3.** Qualitative results on *dVPN* dataset. From left to right, they are left image, right image, right depth map and reconstructed right image.

## 3 Experiments and results

### 3.1 Dataset

We evaluated SADepth on two datasets. The first was the *dVPN* dataset, collected from da Vinci partial nephrectomy, with 34320 pairs of rectified stereo images for training and 14382 pairs for testing [31]. The second was the *SCARED* dataset [1] released during the Endovis challenge at MICCAI 2019, with 17206 pairs (dataset 1, 2, 3, 6 and 7) of rectified stereo images for training and 5637 pairs for testing. To verify the generalization of our framework, we only trained on the *dVPN* dataset but test on both *dVPN* and *SCARED* dataset.

### 3.2 Evaluation Metrics, Baseline, and Implementation Details

**Evaluation Metrics** As the ground truth depth labels were not available for the *in vivo* surgical data in the *dVPN* dataset, we adopted the SSIM index to evaluate the similarity between the reconstructed image and the original input image (*i.e.* $I^{r*}$ and $I^r$) as the evaluation metric. For the *SCARED* dataset the team at Intuitive Surgical collected the ground truth by using structured light, thus we used the absolute error to assess our SADepth model.

**Table 1.** SSIM score for *dVPN* test set (higher is better).

| Method | Training | Mean SSIM | Std. SSIM |
|---|---|---|---|
| ELAS [8] | No training | 47.3 | 0.079 |
| SPS [30] | No training | 54.7 | 0.092 |
| V-Basic [31] | Unsupervised | 55.5 | 0.106 |
| V-Siamese [31] | Unsupervised | 60.4 | 0.066 |
| Monodepth [9] | Unsupervised | 54.9 | 0.087 |
| Monodepth2 [10] | Unsupervised | 71.2 | 0.075 |
| SADepth (ours) | Unsupervised | **79.6** | **0.049** |

**Baseline** We compared SADepth with several recent works. For the *dVPN* dataset, we compared our method with stereo matching-based methods: ELAS [8] and SPS [30]; Siamese-based networks: V-Basic [31] and V-Siamese [31]; and recent deep learning methods: Monodepth [9] and the stereo mode of Monodepth2 [10]. For the *SCARED* dataset, we compared our results with the methods summarized by the recent MICCAI sub-challenge paper [1].

**Implementation Details** The SADepth model was implemented in PyTorch [25], with a batch size of 16 and input/output resolution of $192\times384$. The learning rate was set to $10^{-4}$ for the first 15 epochs and then dropped to $10^{-5}$ for the remainder. The model was trained for 20 epochs using the Adam optimizer which took about 22 hours on a single NVIDIA 2080 Ti GPU.

**Table 2.** The mean absolute depth error for the SCARED test set 1 and 2 (unit: mm) (lower is better).

| Method | Training | Test Set 1 Average | Test Set 2 Average |
|---|---|---|---|
| Lalith Sharan [1] | Supervised | 43.03 | 48.72 |
| Xiaohong Li [1] | Supervised | 22.77 | 20.52 |
| Huoling Luo [1] | Supervised | 19.52 | 18.21 |
| Zhu Zhanshi [1] | Supervised | 9.60 | 21.20 |
| Wenyao Xia [1] | Supervised | 6.73 | 9.44 |
| Congcong Wang [1] | Supervised | 4.10 | 4.28 |
| Trevor Zeffiro [1] | Supervised | 3.60 | 3.47 |
| J.C. Rosenthal [1] | Supervised | 3.44 | 4.05 |
| Dimitris Psychogyios 1 [1] | Supervised | 3.00 | 1.67 |
| Dimitris Psychogyios 2 [1] | Supervised | 2.95 | 2.30 |
| KeXue Fu [1] | Unsupervised | 20.94 | 17.22 |
| Monodepth [9] | Unsupervised | 23.56 | 21.62 |
| Monodepth2 [10] | Unsupervised | 21.92 | 15.25 |
| SADepth (ours) | Unsupervised | **17.42** | **11.23** |

### 3.3   Results

The SADepth and other state-of-the-art results for the *dVPN* dataset are summarized in Table 1 using the mean and standard deviation (Std.) of the SSIM index. The SADepth model effectively outperformed other methods with an SSIM of 79.6, *i.e.* 24.7 units higher than Monodepth [9], 8.4 units higher than Monodepth2 [10], and 19.2 units higher than the Siamese architecture [31].

Table 2 presents the results of SADepth on the test set 1 and test set 2 (as defined in the *SCARED* dataset), together with the performance reported in the MICCAI sub-challenge summary paper [1]. The results show an improvement over the unsupervised methods from the summary paper and recent baselines,

while it is also competitive with some supervised approaches. This confirms that SADepth generalizes well across different datasets collected from different laparoscopes and subjects, while still producing superior performance compared with the state-of-the-art unsupervised approaches.

## 4    Conclusions

We have presented a new self-supervised adversarial depth estimation framework SADepth with an encoder-decoder generator and a concatenated stereo image pair as the input. The adversarial learning strategy improved the generation quality of the framework and led to the state-of-the-art performance on two public datasets. Furthermore, SADepth did not require any per-pixel depth labels and generalized well across different laparoscopes, suggesting excellent applicability to scalable data acquisition when accurate ground truth depth cannot be collected.

## References

1. Allan, M., Mcleod, J., Wang, C.C., Rosenthal, J.C., Fu, K.X., Zeffiro, T., Xia, W., Zhanshi, Z., Luo, H., Zhang, X., et al.: Stereo correspondence and reconstruction of endoscopic data challenge. arXiv:2101.01133 (2021)
2. Chang, J.R., Chen, Y.S.: Pyramid stereo matching network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5410–5418 (2018)
3. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (elus). arXiv preprint arXiv:1511.07289 (2015)
4. Do, T., Nguyen, B.X., Tjiputra, E., Tran, M., Tran, Q.D., Nguyen, A.: Multiple meta-model quantifying for medical visual question answering. arXiv preprint arXiv:2105.08913 (2021)
5. Duggal, S., Wang, S., Ma, W.C., Hu, R., Urtasun, R.: Deeppruner: Learning efficient stereo matching via differentiable patchmatch. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4384–4393 (2019)
6. Eigen, D., Puhrsch, C., Fergus, R.: Depth map prediction from a single image using a multi-scale deep network. arXiv preprint arXiv:1406.2283 (2014)
7. Garg, R., Bg, V.K., Carneiro, G., Reid, I.: Unsupervised cnn for single view depth estimation: Geometry to the rescue. In: European conference on computer vision. pp. 740–756. Springer (2016)
8. Geiger, A., Roser, M., Urtasun, R.: Efficient large-scale stereo matching. In: Asian conference on computer vision. pp. 25–38. Springer (2010)
9. Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 270–279 (2017)
10. Godard, C., Mac Aodha, O., Firman, M., Brostow, G.J.: Digging into self-supervised monocular depth estimation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3828–3838 (2019)
11. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks. arXiv preprint arXiv:1406.2661 (2014)

12. Grasa, O.G., Bernal, E., Casado, S., Gil, I., Montiel, J.: Visual slam for hand-held monocular endoscope. IEEE transactions on medical imaging **33**(1), 135–146 (2013)
13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
14. Heise, P., Klose, S., Jensen, B., Knoll, A.: Pm-huber: Patchmatch with huber regularization for stereo matching. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2360–2367 (2013)
15. Huang, B., Tsai, Y.Y., Cartucho, J., Vyas, K., Tuch, D., Giannarou, S., Elson, D.S.: Tracking and visualization of the sensing area for a tethered laparoscopic gamma probe. International Journal of Computer Assisted Radiology and Surgery **15**(8), 1389–1397 (2020)
16. Huang, B., Zheng, J.Q., Giannarou, S., Elson, D.S.: H-net: Unsupervised attention-based stereo depth estimation leveraging epipolar geometry. arXiv preprint arXiv:2104.11288 (2021)
17. Jaderberg, M., Simonyan, K., Zisserman, A., Kavukcuoglu, K.: Spatial transformer networks. arXiv preprint arXiv:1506.02025 (2015)
18. Johnston, A., Carneiro, G.: Self-supervised monocular trained depth estimation using self-attention and discrete disparity volume. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4756–4765 (2020)
19. Joung, S., Kim, S., Park, K., Sohn, K.: Unsupervised stereo matching using confidential correspondence consistency. IEEE Transactions on Intelligent Transportation Systems **21**(5), 2190–2203 (2019)
20. Leonard, S., Sinha, A., Reiter, A., Ishii, M., Gallia, G.L., Taylor, R.H., Hager, G.D.: Evaluation and stability analysis of video-based navigation system for functional endoscopic sinus surgery on in vivo clinical data. IEEE transactions on medical imaging **37**(10), 2185–2195 (2018)
21. Liu, X., Sinha, A., Ishii, M., Hager, G.D., Reiter, A., Taylor, R.H., Unberath, M.: Dense depth estimation in monocular endoscopy with self-supervised learning methods. IEEE transactions on medical imaging **39**(5), 1438–1447 (2019)
22. Luo, W., Schwing, A.G., Urtasun, R.: Efficient deep learning for stereo matching. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5695–5703 (2016)
23. Mack, M.J.: Minimally invasive and robotic surgery. Jama **285**(5), 568–572 (2001)
24. Nguyen, A., Kundrat, D., Dagnino, G., Chi, W., Abdelaziz, M.E., Guo, Y., Ma, Y., Kwok, T.M., Riga, C., Yang, G.Z.: End-to-end real-time catheter segmentation with optical flow-guided warping during endovascular intervention. In: 2020 IEEE International Conference on Robotics and Automation (ICRA). pp. 9967–9973. IEEE (2020)
25. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch (2017)
26. Pilzer, A., Xu, D., Puscas, M., Ricci, E., Sebe, N.: Unsupervised adversarial depth estimation using cycled generative networks. In: 2018 International Conference on 3D Vision (3DV). pp. 587–595. IEEE (2018)
27. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. pp. 234–241. Springer (2015)
28. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. IEEE transactions on image processing **13**(4), 600–612 (2004)

29. Watson, J., Firman, M., Brostow, G.J., Turmukhambetov, D.: Self-supervised monocular depth hints. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2162–2171 (2019)
30. Yamaguchi, K., McAllester, D., Urtasun, R.: Efficient joint segmentation, occlusion labeling, stereo and flow estimation. In: European Conference on Computer Vision. pp. 756–771. Springer (2014)
31. Ye, M., Johns, E., Handa, A., Zhang, L., Pratt, P., Yang, G.Z.: Self-supervised siamese learning on stereo image pairs for depth estimation in robotic surgery. arXiv preprint arXiv:1705.08260 (2017)
32. Yi, Z., Zhang, H., Tan, P., Gong, M.: Dualgan: Unsupervised dual learning for image-to-image translation. In: Proceedings of the IEEE international conference on computer vision. pp. 2849–2857 (2017)
33. Zhang, D., Xiao, B., Huang, B., Zhang, L., Liu, J., Yang, G.Z.: A self-adaptive motion scaling framework for surgical robot remote control. IEEE Robotics and Automation Letters **4**(2), 359–366 (2018)
34. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1851–1858 (2017)
35. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: Proceedings of the IEEE international conference on computer vision. pp. 2223–2232 (2017)