

# Camera Augmentation: Enabling Uncalibrated Stereo Matching of Minimally-Invasive Surgery Images by Training from the Wealth of Public Synthetic Image Datasets

Rasoul Sharifian<sup>1,2\*†</sup>, Navid Rabbani<sup>1†</sup>, Yongcong Zhang<sup>1,2</sup>,  
Adrien Bartoli<sup>1,2</sup>

<sup>1</sup>\*DIA2M, DRCI, CHU Clermont-Ferrand, France.

<sup>2</sup>SURGAR, Surgical Augmented Reality, Clermont-Ferrand, France.

\*Corresponding author(s). E-mail(s): [rasoul.sharifian.cs@gmail.com](mailto:rasoul.sharifian.cs@gmail.com);

†These authors contributed equally to this work.

## Abstract

**Purpose.** Existing models for stereo matching in Minimally-Invasive Surgery (MIS) require calibrated stereo images. Accurate calibration is however often unavailable intraoperatively. Training an uncalibrated stereo model is thus attractive but challenging owing to the lack of disparity-labelled surgical images.

**Methods.** We leverage the wealth of non-medical stereo synthetic image datasets. These data were however generated in ideal conditions –rectified and with centred principal points– hence differ from real uncalibrated MIS images. We propose camera augmentation, a new type of image augmentation that augments a dataset by altering the camera’s orientation and intrinsic parameters via geometric parameters. We augment the idealised existing datasets, sampling the geometric augmentation parameters from distributions estimated through an in-depth analysis and modelling of stereo laparoscopes. This forms the Camera Augmentation Training Strategy (CATS), with which we retrain RAFTStereo and IGEV++ for zero-shot uncalibrated stereo matching in MIS.

**Results.** We evaluated using the SCARED, StereoMIS, RIS2017, and an in-house datasets. In the uncalibrated setting on the SCARED dataset, CATS-RAFTStereo and CATS-IGEV++ achieved End-Point-Errors (EPE) of 1.42 and 1.41 pixels. This is a successful result, as the reference pretrained models obtained 1.23 and 1.21 pixels in the calibrated setting and failed in the uncalibrated setting.

001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046

047 **Conclusion.** Camera augmentation bridges the gap between ideally-conditioned  
048 datasets and the real surgical conditions of uncertain or unavailable calibration,  
049 enabling the retraining of state-of-the-art architectures. Beyond stereo, the pro-  
050 posed CATS is applicable to tasks sensitive to camera geometry. Code and models  
051 will be released publicly.

052 **Keywords:** Minimally Invasive Surgery, Stereo Matching, Disparity, Uncalibrated  
053

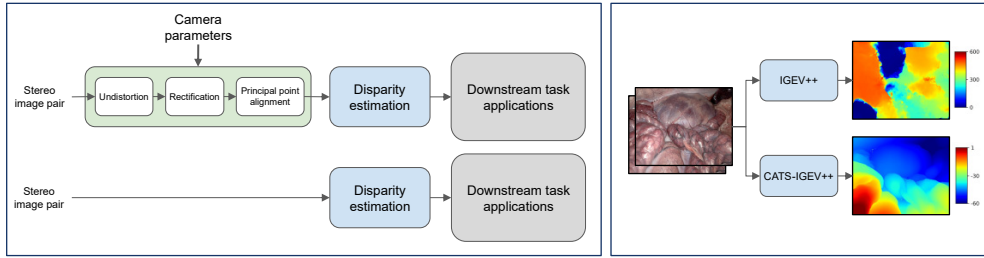
054  
055

## 056 1 Introduction

057  
058 Depth estimation in MIS plays a crucial role in downstream applications such as Aug-  
059 mented Reality [1–3] and robotic surgery guidance [3]. Even if depth can be estimated  
060 from monocular images acquired by standard laparoscopes, stereo depth estimation  
061 benefits from strong triangulation constraints arising from left–right image consistency.  
062 This reduces the inaccuracies commonly occurring in monocular depth estimation  
063 where the models heavily rely on learnt priors. The use of stereo laparoscopes has  
064 been steadily increasing, owing to the rise of robot-assisted MIS systems.

065 A stereo system comprises two cameras. The stereo calibration process determines  
066 each camera’s intrinsic parameters and the extrinsic parameters describing their rela-  
067 tive rotation and translation. Calibrated stereo methods have standard preprocessing  
068 steps: undistortion, rectification, and principal-point alignment (Fig. 1). Disparity esti-  
069 mation is then performed, and the disparity mapped to metric depth. While early  
070 works established the fundamental components [4], recent neural architectures such  
071 as RAFTStereo [5], IGEV++ [6] and FoundationStereo [7] have shown substantially  
072 improved performance and robustness and, importantly, generalisation across datasets.  
073 These models require rectified images with centred principal points, which necessitate  
074 camera calibration; as expected, their performance directly depends on the accuracy  
075 of calibration parameters. Calibrated stereo methods specifically adapted to MIS were  
076 proposed, including MSDESIS [8] and MCF-SMSIS [9] which jointly estimate dispar-  
077 ity and segmentation with shared features. Beyond the medical domain, recent stereo  
078 matching methods have demonstrated promising cross-domain generalisation [5–7]. In  
079 practice, however, accurate calibration may not be available throughout surgery, as  
080 the parameters may vary with zoom level, focusing and any other digital or mechan-  
081 ical adjustments. Calibration is also impractical to perform routinely in operating  
082 room conditions. Consequently, both the intrinsic and extrinsic parameters are often  
083 unavailable or unreliable, making the standard stereo pipeline inapplicable. In spite of  
084 these known practical challenges in camera calibration, uncalibrated stereo depth esti-  
085 mation –independent of camera calibration parameters, both intrinsics and extrinsics–  
086 remains largely unexplored in MIS. We show that, while differences in scene content  
087 may not pose a major issue for large learning-based models, discrepancies in cam-  
088 era settings ultimately lead to systematic failures. This issue cannot be resolved by  
089 fine-tuning over the limited depth-labelled MIS data.

090 In the uncalibrated setting (Fig. 1), although the resulting disparity cannot be  
091 mapped to metric depth, as this mapping depends on the calibration parameters, it  
092



**Fig. 1** Left: Pipelines for calibrated (top) and uncalibrated (bottom) stereo matching. Right: disparity estimated by the IGEV++ [6] and the proposed CATS-IGEV++ models in uncalibrated setting.

can still be mapped to *relative* depth. Importantly, many downstream tasks in MIS operate effectively with relative rather than metric depth (*e.g.* collision avoidance [10], instrument tracking [11] and surgical simulation [12]). Augmented reality systems compute the scale of uncalibrated reconstructions concurrently with registration, and visual servoing controls robots effectively from uncalibrated images and reconstructions, highlighting the practical utility of relative depth. We address the problem of uncalibrated stereo in MIS by leveraging the wealth of non-medical stereo synthetic image datasets. A major challenge is that these data were generated in ideal conditions (perfectly rectified, with centred principal point) and published without 3D scene parameters, preventing re-rendering. To overcome these limitations, we introduce the Camera Augmentation Training Strategy (CATS), a framework that enables robust stereo learning under uncalibrated conditions. Concretely, our main contributions are threefold. First, we propose *camera augmentation*, a framework that augments a dataset by altering the camera’s orientation and intrinsic parameters. It does not require knowing the camera pose or the depth; yet, it offers geometric control over the augmented data. Second, we use camera augmentation to enhance existing idealised non-medical datasets to simulate uncalibrated yet near-rectified stereo laparoscopes. We sample geometric parameters from distributions derived through a detailed analysis and modelling of stereo laparoscope configurations. Third, we use CATS to retrain RAFTStereo and IGEV++ for zero-shot uncalibrated stereo matching in MIS.

## 2 Camera Augmentation

We introduce camera augmentation as a generalisation of conventional image augmentation, in which the augmented images are generated by perturbing the camera parameters. Unlike conventional augmentation performed by basic 2D geometric transformations within the image plane, camera augmentation explicitly models the geometric image-formation process. This allows the generation of new, yet geometrically-consistent image samples, thereby enriching datasets that lack sufficient diversity in camera orientation, focal length, principal point, or distortion.

Generating images is challenging, as it typically relies on rendering, which requires at least the lighting, 3D scene surface, reflection function, and camera parameters. This complex process is incompatible with the typical practical conditions in which most

139 of these parameters are unavailable. In camera augmentation, we bypass this complex  
 140 procedure by simply keeping the original camera position unchanged, which nonethe-  
 141 less leaves a variety of other geometric controls, subsuming conventional augmentation.  
 142 These new controls allow the generation of images that preserve 3D-geometric con-  
 143 sistency without requiring complex parameters and computations, yet still improve  
 144 model robustness to changes in camera geometry.

145 We model the original camera using its intrinsic parameters  $\mathbf{K} \in \mathbb{R}^{3 \times 3}$  and its  
 146 distortion parameters  $\boldsymbol{\theta} \in \mathbb{R}^k$ . Choosing as 3D coordinate frame  $\mathcal{C}$  the standard coor-  
 147 dinate frame related to this camera (origin centred on the projection centre,  $Z$  axis  
 148 coincident with the optical axis, and  $X, Y$  axes parallel to the retinal plane axes), the  
 149 camera pose is the identity rotation  $\mathbf{I}$  and zero translation  $\mathbf{0}$ . We model the augmented  
 150 image camera using its intrinsics  $\mathbf{K}' \in \mathbb{R}^{3 \times 3}$ , its distortion parameters  $\boldsymbol{\theta}' \in \mathbb{R}^k$  and its  
 151 pose in  $\mathcal{C}$  as the rotation  $\mathbf{R}$ . We leave, as already stated, the translation to  $\mathbf{0}$  in order  
 152 to keep the camera position unchanged.

153 Without distortion, for  $\mathbf{P}$  the Euclidean coordinates of a 3D point, and  $\tilde{\mathbf{P}}$  its homo-  
 154 geneous coordinates, the projections are  $\tilde{\mathbf{p}} \sim \mathbf{K} [\mathbf{I} \ \mathbf{0}] \tilde{\mathbf{P}} \sim \mathbf{K}\mathbf{P}$  and  $\tilde{\mathbf{p}}' \sim \mathbf{K}' [\mathbf{R} \ \mathbf{0}] \tilde{\mathbf{P}} \sim$   
 155  $\mathbf{K}'\mathbf{R}\mathbf{P}$ , where ' $\sim$ ' means equality up to scale. By inverting the first projection and  
 156 substituting into the second one, we obtain the well-known homographic relationship  
 157  $\tilde{\mathbf{p}}' \sim \mathbf{H}\tilde{\mathbf{p}}$ , with  $\mathbf{H} \sim \mathbf{K}'\mathbf{R}\mathbf{K}^{-1}$ . We can factor the homography as  $\mathbf{H} = \mathbf{K}_{\text{rel}}\mathbf{H}_{\text{R}}$  where  
 158  $\mathbf{K}_{\text{rel}} = \mathbf{K}'\mathbf{K}^{-1}$  and  $\mathbf{R}$  represent the sought *relative intrinsic* and *relative rotation*, and  
 159  $\mathbf{H}_{\text{R}} = \mathbf{K}\mathbf{R}\mathbf{K}^{-1}$ . Using the focal length ratios  $s_x = f'_x/f_x$  and  $s_y = f'_y/f_y$ , we have:

$$160$$

$$161 \quad \mathbf{K} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{K}' = \begin{bmatrix} f'_x & 0 & c'_x \\ 0 & f'_y & c'_y \\ 0 & 0 & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{K}_{\text{rel}} = \begin{bmatrix} s_x & 0 & c'_x - s_x c_x \\ 0 & s_y & c'_y - s_y c_y \\ 0 & 0 & 1 \end{bmatrix}.$$

$$162$$

$$163$$

164 Therefore,  $\mathbf{K}_{\text{rel}}$  is independent of the absolute focal lengths. Conventional image aug-  
 165 mentations, including random rotation, scaling, and translation within the image  
 166 plane, can be interpreted as special cases of the proposed camera augmentation.

167 We now add lens distortion, for which we define the generic distortion and  
 168 undistortion function  $\mathcal{D}_{\mathbf{K},\boldsymbol{\theta}}(\cdot)$  and  $\mathcal{D}_{\mathbf{K},\boldsymbol{\theta}}^{-1}(\cdot)$ . These functions depend on the intrinsic  
 169 parameters in  $\mathbf{K}$  and the distortion coefficients  $\boldsymbol{\theta}$ , which typically model radial and tan-  
 170 gential distortions. For the perspective function  $\Pi([q_1, q_2, q_3]^\top) = [q_1/q_3, q_2/q_3]^\top$  and  
 171 the homogenisation function  $\Gamma([q_1, q_2]^\top) = [q_1, q_2, 1]^\top$ , the complete transformation is:

$$172$$

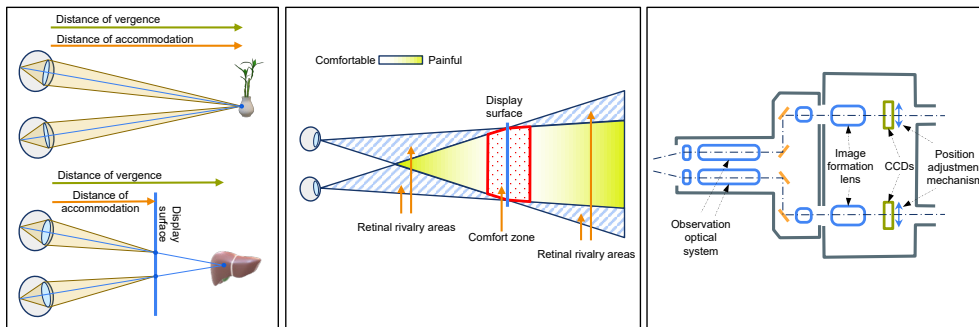
$$173 \quad \mathbf{p}' \sim \mathcal{D}_{\mathbf{K}',\boldsymbol{\theta}'} \left( \Pi \left( \mathbf{K}_{\text{rel}} \mathbf{H}_{\text{R}} \Gamma \left( \mathcal{D}_{\mathbf{K},\boldsymbol{\theta}}^{-1}(\mathbf{p}) \right) \right) \right). \quad (1)$$

$$174$$

175 In performing camera augmentation,  $\mathbf{K}_{\text{rel}}$ ,  $\mathbf{R}$  and  $\boldsymbol{\theta}'$  are drawn from probabilistic dis-  
 176 tributions that model realistic variations in camera parameters in the target domain.  
 177 These distributions act as a bridge between the known camera configuration in the  
 178 available dataset and the perturbed configuration in the augmented setup, allow-  
 179 ing camera augmentation to be performed without re-rendering. The use with stereo  
 180 laparoscopes is discussed in the next sections, where we measure that lens distortion  
 181 is negligible for the public datasets at hand. We thus use the following *distortion-free*  
 182 *simplification* of Eq. 1:

$$183 \quad \mathbf{p}' \sim \mathbf{K}_{\text{rel}} \mathbf{H}_{\text{R}} \mathbf{p}. \quad (2)$$

$$184$$



**Fig. 2** Left: Vergence–accommodation conflict (VAC); In natural viewing, vergence and accommodation coincide (top). In stereo displays, vergence targets the virtual object while accommodation stays on the display plane (bottom). Middle: Combined effects of binocular rivalry and vergence–accommodation conflict, highlighting the range where both are minimised. Right: CCD position adjustment mechanism in a stereo laparoscope bringing the virtual image into the visual comfort zone.

### 3 Perception and the Design of Stereo Laparoscopes

We analyse the perceptual constraints and the design of stereo laparoscopes. This establishes a basis for modelling variations in intrinsic and extrinsic parameters. The design choices of stereo laparoscopes given below are general and may not be universally followed by all manufacturers.

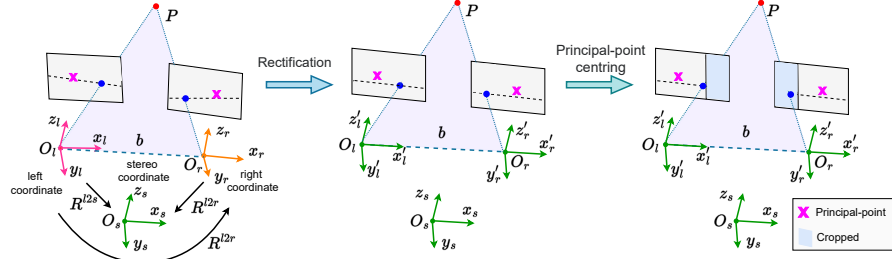
#### *Overview*

Depth perception in humans may be enabled by displaying stereo images, should some perceptual constraints be satisfied. A stereo laparoscope is a dual-channel imaging system that captures synchronised image pairs for such displays. Its base components, shown in Fig. 2, are two parallel optical pathways with relay lenses, fibre-based illumination, and twin sensors. The sensors are factory-calibrated to provide left and right views for depth perception via a binocular eyepiece or 3D display. The surgeon’s visual discomfort, including fatigue and eye strain, and quality of depth perception, directly depend on the extent to which the stereo imaging pipeline’s design satisfies the perceptual constraints [13]. These constraints and design choices can be gathered in two groups, which we describe next, along with an experimental verification.

#### *Geometric Alignment*

Depth perception requires precise alignment of the left and right images. In an ideal stereo configuration, correspondences exhibit purely horizontal displacements. In computer vision, this alignment is also referred to as rectified images, whose epipolar lines are horizontal. The careful design and factory-based placement adjustment of the imaging lenses and sensors allow modern stereoscopic laparoscopes to produce near-rectified image pairs. We quantified the residual misalignment experimentally from the camera parameters given in the public datasets SCARED [14], StereoMIS [15] and RIS2017 [16], which were acquired using da Vinci systems. We used the extrinsic parameters to estimate three rotations around each camera’s optical centre. These

185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230



**Fig. 3** Geometry of stereo camera pairs in unrectified (left), rectified (middle), and principal-point aligned (right) configurations. The principal point shift and scan-line misalignment are intentionally exaggerated for visualisation purposes. In practice, the applied perturbations are much smaller and preserve sufficient overlap between the two views.

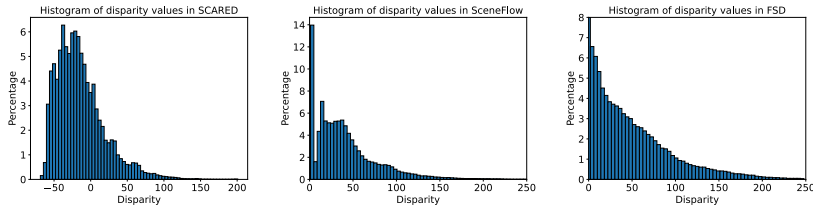
rotations are computed relative to the stereo coordinate system shown in Fig. 3, defined with the  $X$ -axis aligned with the baseline and taken as minimal relative to the left camera. Statistics -mean and standard deviation- for the sampled angles are given in Table 1. Their concentration around zero, all being much lower than 1 degree, and similarity across datasets indicate that different laparoscopes share a comparable small deviation from the rectified design. As indicated by the small  $Y$ -axis rotations in Table 1, the display planes follow a planar design without toeing-in. We used the intrinsic parameters to quantify lens distortion, which is known to impair rectification. We measured a residual distortion consistently below half a pixel; given the image resolution of  $1280 \times 1024$ , this confirms a near distortion-free design. Lastly, the principal point for both cameras must be aligned; this is discussed below.

### *Vergence-Accommodation Conflict and Binocular Rivalry*

Vergence rotates the eyes to enable single vision, while accommodation adjusts focus. In natural conditions, both occur at the same distance. In stereoscopic displays however, VAC arises from vergence targetting the virtual object while accommodation remains fixed on the display plane. Binocular rivalry occurs when the left and right images contain inconsistent information, often due to asymmetric fields of view. Its extent increases as the virtual object approaches the viewer. Fig. 2 shows a depth range where the effects of both VAC and binocular rivalry are minimised. Virtual anatomical structures are typically rendered slightly in front of or behind the display plane to optimise depth perception and comfort. The way these two perceptual constraints are coped with in stereo laparoscopes is by shifting the sensors during manufacturer calibration, as documented in patented systems [17–19]. This shift is chosen to place maximal comfort at the typical surgical depth of 5 cm. In terms of the camera parameters, the effect is to offset the principal points horizontally away from the image centre. This creates a camera-induced domain gap between stereo laparoscope images and conventional stereo images, where principal points are centred. A direct effect is that the stereo disparity may be positive and negative. In contrast, state-of-the-art stereo matching methods are typically trained on datasets with only positive disparities, such as SceneFlow [20] and FSD [7], whereas real laparoscopic datasets like SCARED [14] contain both positive and negative disparities. Fig. 4 illustrates this gap.

**Table 1** Statistics on camera pose deviations from the ideal rectified configuration for the left (top row) and right (bottom row) cameras, for the three Euler rotation angles, in degrees.

	$R_x^{l2s}$		$R_y^{l2s}$		$R_z^{l2s}$		$R_x^{r2s}$		$R_y^{r2s}$		$R_z^{r2s}$	
	$\mu$	$\sigma$	$\mu$	$\sigma$	$\mu$	$\sigma$	$\mu$	$\sigma$	$\mu$	$\sigma$	$\mu$	$\sigma$
SCARED	0.0000	0.0001	-0.0325	0.0129	0.0904	0.3597	-0.0014	0.0053	0.0013	0.0266	0.0910	0.3602
StereoMIS	-0.0006	0.0006	0.3673	0.0032	-0.1907	0.1812	-0.0020	0.0008	0.3799	0.0187	-0.1887	0.1809
RIS2017	-0.0004	0.0014	-0.3386	0.1652	0.2637	0.4680	0.0002	0.0076	-0.2744	0.1640	0.2639	0.4681
SCARED+StMIS	-0.0001	0.0004	0.0564	0.1666	0.0280	0.3487	-0.0016	0.0047	0.0854	0.1594	0.0289	0.3489



**Fig. 4** Disparity distributions (in px) for the SCARED, SceneFlow, and FSD datasets.

## 4 Camera Augmentation in Stereo Laparoscopes

In classical stereo vision, disparity is defined for rectified image pairs. Surgical laparoscopes produce near-rectified images, as described in Section 3, where vertical displacements are small and horizontal displacements can be used as estimates for disparity. Disparity is related to depth through  $d = \frac{fb}{z} + (c_x^L - c_x^R)$ , where  $f$  is the focal length,  $b$  is the baseline, and  $c_x^L$  and  $c_x^R$  are the  $x$ -coordinates of the left and right principal points. This relationship shows that depth in uncalibrated settings is only available *relatively*, up to an unknown scale and shift.

We use camera augmentation from Section 2 to address uncalibrated stereo matching in MIS. The proposed method, termed Camera-Augmentation Training Strategy (CATS), is used to retrain existing stereo matching architectures on camera-augmented datasets that simulate realistic imaging imperfections. We retrained RAFTStereo and IGEV++ using the original training protocol, hyperparameters, and datasets, with as sole modification the use of CATS. Both models are trained on the SceneFlow dataset [20], which contains over 39K stereo pairs with a resolution of  $960 \times 540$  pixels, rendered from a wide range of 3D scenes. The SceneFlow dataset uses ideal stereo camera geometry, including perfect rectification, centred principal points, and the absence of lens distortion, differing markedly from the real MIS conditions.

We use the analysis of MIS stereo camera imperfections from Section 3 (deviations in rectification and principal point misalignment), to generate camera augmentations. This augmentation strategy encourages the network to learn a stereo matching process adapted to uncalibrated MIS images. Concretely, for each of the left and right cameras, we sample three Euler angles,  $\phi$ ,  $\theta$ , and  $\psi$ , from Gaussian distributions reflecting the

323 measured perturbation statistics, shown in Table 1:

324

$$325 \quad \phi_c \sim \mathcal{N}(\mu_{c,\phi}, \sigma_{c,\phi}^2), \quad \theta_c \sim \mathcal{N}(\mu_{c,\theta}, \sigma_{c,\theta}^2), \quad \psi_c \sim \mathcal{N}(\mu_{c,\psi}, \sigma_{c,\psi}^2), \quad c \in \{\text{left}, \text{right}\}. \quad (3)$$

326

327 The rotation matrix  $\mathbf{R}$  is constructed from the Euler angles  $(\phi, \theta, \psi)$  following the  $Z, Y,$   
 328  $X$  (yaw-pitch-roll) convention as  $\mathbf{R}_c = \mathbf{R}_z(\psi_c) \mathbf{R}_y(\theta_c) \mathbf{R}_x(\phi_c)$ ,  $c \in \{\text{left}, \text{right}\}$ . The  
 329 corresponding homography induced by the rotational perturbation for each camera is  
 330 computed as:

331

$$331 \quad \mathbf{H}_{\mathbf{R},c} = \mathbf{K} \mathbf{R}_c \mathbf{K}^{-1}, \quad \forall c \in \{\text{left}, \text{right}\}, \quad (4)$$

332

333 where  $\mathbf{K}$  is the intrinsic calibration matrix from the original dataset. We then sample  
 334  $\mathbf{K}_{\text{rel}}$  for the left and right cameras. Since random scaling is already included in the  
 335 standard data augmentation pipelines of the reference methods, additional variation  
 336 in scaling factors,  $s_x$  and  $s_y$ , would be redundant, and we use  $s_x = s_y = 1$  for both  
 337 cameras. Furthermore, based on observations from Section 3, we use  $c'_y = c_y$  for both  
 338 cameras, as the vertical principal point displacement was found to be negligible. The  
 339 dominant parameter affecting  $\mathbf{K}_{\text{rel}}$  is the horizontal offset between the principal points  
 340 of the left and right cameras. Accordingly, we set  $c'_{x,\text{left}} = c_{x,\text{left}}$  and sample the right  
 341 camera offset,  $\Delta x_{\text{right}} = c'_{x,\text{right}} - c_{x,\text{right}}$  from a uniform distribution estimated from  
 342 the empirical distribution given in Fig. 4 as  $\Delta x_{\text{right}} \sim \mathcal{U}(-100, 0)$ . The uniform dis-  
 343 tribution is a pragmatic choice ensuring coverage of the observed empirical range and  
 344 keeping the sampling simple. The use of non-parametric distributions would greatly  
 345 increase complexity. This distribution strikes a balance between positive and negative  
 346 disparities during training, improving robustness across varying stereo configurations.  
 347 We arrive at the following sampled relative intrinsic matrices:

348

349

$$349 \quad \mathbf{K}_{\text{rel},\text{left}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{K}_{\text{rel},\text{right}} = \begin{bmatrix} 1 & 0 & \Delta x_{\text{right}} \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (5)$$

350

351 Finally, according to observations from Section 3, the effects of lens distortion are  
 352 negligible within the analysed MIS datasets; hence, distortion is omitted in CATS.  
 353 Importantly, we apply the same homography as for the left image to the disparity  
 354 ground truth and then adjust all disparity values by adding  $\Delta x_{\text{right}}$ .

355

## 356 5 Experiments and Results

357

358 **Reference models.** We selected three stereo matching models representing the SoTA:  
 359 RAFTStereo [5], IGEV++ [6], and FoundationStereo [7] (we used the best performing  
 360 version for general use, based on ViT-Large), which we used as released, without addi-  
 361 tional training or fine-tuning. These architectures represent the leading advances in  
 362 general-purpose stereo matching, independently of our extending them to uncalibrated  
 363 MIS. For comparison with models specifically designed for MIS, we included MSDE-  
 364 SIS [8] and MCF-SMSIS [9], representing the SoTA in MIS stereo matching. We used  
 365 RAFT [21] as a representative optical flow method, taking the vertical component of  
 366 the predicted optical flow as an estimate of disparity. Finally, we used DUST3R [22]  
 367 as a representative general multi-view reconstruction method.

368

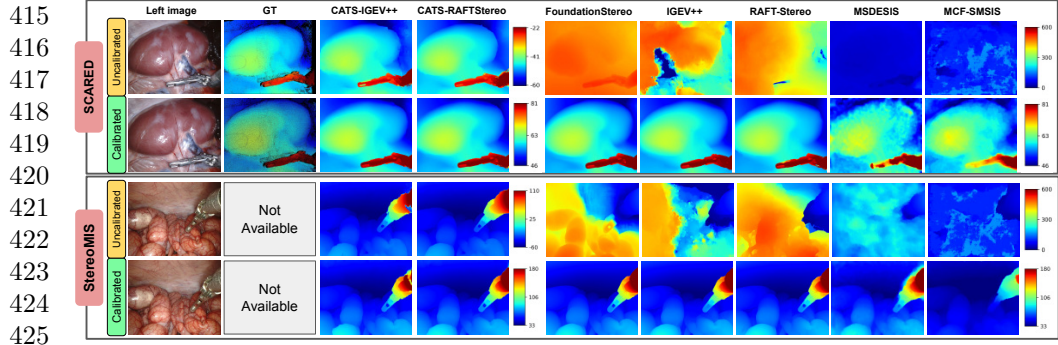
**Table 2** Uncalibrated scenario evaluation over original images from the SCARED dataset. Depth values are computed from predicted disparity using static calibration parameters for reference only. In the DUS3R method, the calibration parameters are used to find disparity values from the predicted depth.

	Disparity metrics (px)		Depth metrics (mm)			
	EPE ↓	Bad 3.0 ↓	MAE ↓	RMSE ↓	AbsRel ↓	$\delta_{0.25}$ ↑
RAFTStereo	244.37	80.08	43.95	48.69	0.55	0.23
IGEV++	186.49	79.79	39.75	46.03	<u>0.49</u>	<u>0.24</u>
FoundationStereo	331.41	83.60	48.50	52.05	0.62	0.20
MSDESI	88.25	99.96	<u>38.06</u>	40.42	0.53	0.001
MCF-SMSIS	119.01	97.78	43.38	46.51	0.59	0.03
DUS3R (GT calib.)	67.81	100	232.78	233.30	4.17	0.00
DUS3R (DUS3R calib.)	687.97	100	232.78	233.30	4.17	0.00
RAFT	3.00	27.79	2.85	4.98	0.039	0.81
CATS-RAFTStereo	<u>1.42</u>	<u>6.36</u>	<b>1.33</b>	<b>2.25</b>	<b>0.01</b>	<b>0.97</b>
CATS-IGEV++	<u>1.41</u>	<b>6.32</b>	<b>1.33</b>	<u>2.31</u>	<b>0.01</b>	<b>0.97</b>

**CATS models.** We retrained RAFTStereo and IGEV++ using the proposed CATS, yielding the CATS-RAFTStereo and CATS-IGEV++ models. We did not retrain FoundationStereo as its training code is not yet publicly available; however CATS would be fully compatible with its pipeline since it operates purely at the data level.

**Experimental setup.** We evaluated quantitatively on the SCARED dataset and qualitatively on the SCARED, StereoMIS and an in-house datasets. The in-house dataset was collected from hepatectomy procedures with ethical approval IRB00008526-2019-CE58 issued by CPP Sud-Est VI in Clermont-Ferrand, France. In quantitative evaluation, we used keyframes of the SCARED dataset, excluding video sequences that exhibit inconsistencies between the visual data and the kinematic ground truth [14]. In addition, we excluded datasets 4 and 5, as they are explicitly reported in the SCARED paper to contain errors [14]. In the calibrated scenario, we used the available intrinsic and extrinsic parameters to rectify the stereo pairs and centre the principal points before applying the stereo matching models. In the uncalibrated scenario, we applied stereo matching directly on the original stereo image pairs. We evaluated quantitatively using the End-Point Error (EPE) and Bad-3.0 as disparity metrics, and MAE, RMSE, AbsRel, and  $\delta_{0.25}$  as depth metrics [23].

**Results.** Quantitative results for the uncalibrated and calibrated scenarios are given in Tables 2 and 3. Several key observations can be made. First, as shown in Table 2, all reference and MIS-specific models exhibit systematic failures in the uncalibrated setting, with large disparity and depth errors. RAFTStereo, IGEV++, and FoundationStereo yield EPEs exceeding 180 pixels, and MIS-specific models, MSDESI and MCF-SMSIS, perform only marginally better with EPEs of 88.25 pixels and 119.01 pixels, respectively, still considered failures. DUS3R exhibits an EPE of 67.81 pixels when using the camera parameters provided by the dataset, and 687.97 pixels when using its own estimated parameters, both of which are far from acceptable. RAFT achieves a substantially lower EPE of 3.00 pixels, which is nonetheless more than twice as large as the CATS methods. The depth metrics confirm the failures, with the lowest MAE value



**Fig. 5** Qualitative model comparison in the uncalibrated and calibrated scenarios on samples from the SCARED and StereoMIS datasets.

at 38 mm, indicating poor geometric consistency. Second, the CATS-trained models perform successfully on the same uncalibrated scenario, achieving EPEs of 1.42 pixels for CATS-RAFTStereo and 1.41 pixels for CATS-IGEV++. Depth metrics further confirm this improvement with an MAE of about 1.33 mm and an AbsRel of about 0.01, demonstrating that CATS successfully enables reliable stereo matching directly from uncalibrated MIS stereo pairs. Third, the performance of CATS-trained models in the uncalibrated scenario (EPEs of 1.42 and 1.41 pixels for CATS-RAFTStereo and CATS-IGEV++), when compared to the best-performing method in the calibrated setting in Table 3 (EPE of 1.17 pixels for FoundationStereo), underscores the effectiveness of CATS. Remarkably, CATS-trained models achieve nearly equivalent performance despite the absence of camera calibration. Fourth, under the calibrated setting, the same models achieve comparable performance. CATS-RAFTStereo and CATS-IGEV++ reach EPEs of 1.23 and 1.21 pixels, which are on par with their reference counterparts, with RAFTStereo at 1.21 pixels and IGEV++ at 1.23 pixels. This shows that CATS does not compromise performance when calibration is available. Fifth, in the calibrated scenario, CATS-RAFTStereo and CATS-IGEV++ perform nearly as well as FoundationStereo, the best-performing model, with EPEs marginally higher by 0.04 and 0.06 pixels.

Qualitative results in Fig. 5 and 6 illustrate these findings. In the uncalibrated scenario, reference and MIS-specific models produced distorted and inconsistent disparity maps, whereas the proposed CATS models successfully recovered fine structural details and coherent patterns. In the calibrated scenario, all models performed well, including the CATS variants, confirming that CATS remains fully compatible with calibrated conditions and does not degrade performance. Furthermore, the CATS models demonstrated consistent performance across all tested datasets acquired with different da Vinci system models, highlighting their generalisability. Beyond its strong performance on the three public datasets collected on animal models, CATS produced high-quality, structurally coherent disparity maps on the in-house data of human hepatectomy procedures. This demonstrates the method’s robustness and translational potential, supporting its use in real-world clinical applications such as computer-assisted interventions and intraoperative guidance.

**Table 3** Calibrated scenario evaluation on undistorted, rectified, and principal point-aligned images from the SCARED dataset.

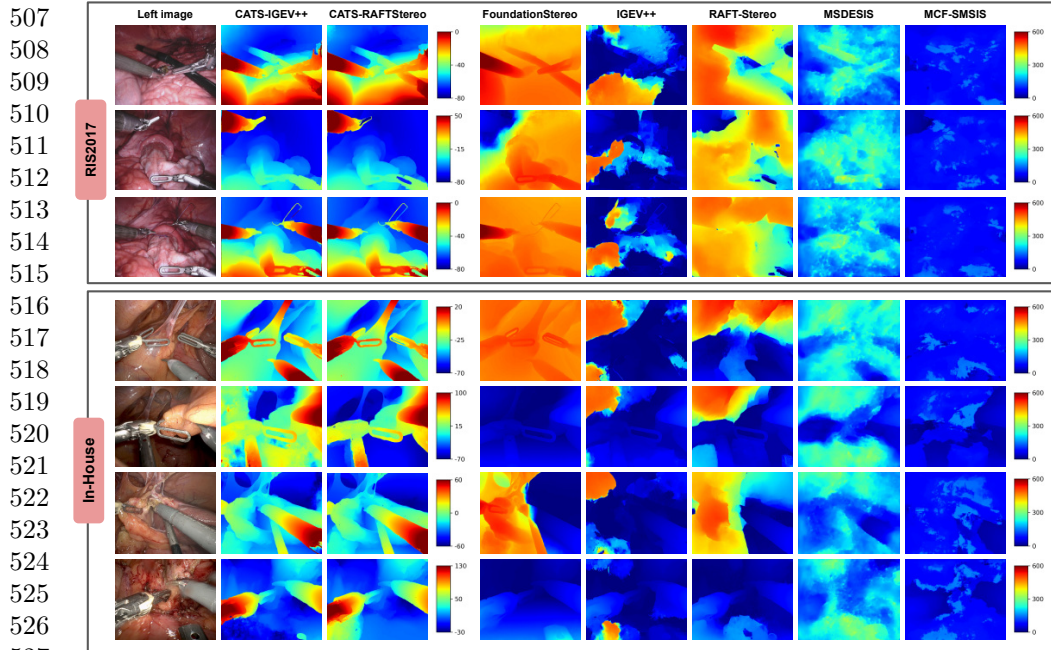
	Disparity metrics (px)		Depth metrics (mm)			
	EPE ↓	Bad 3.0 ↓	MAE ↓	RMSE ↓	AbsRel ↓	$\delta_{0.25}$ ↑
RAFTStereo	<u>1.21</u>	5.67	<u>1.01</u>	<u>2.32</u>	<b>0.014</b>	0.979
IGEV++	1.23	5.56	1.03	2.43	<u>0.015</u>	<u>0.980</u>
FoundationStereo	<b>1.17</b>	<b>5.10</b>	<b>0.97</b>	<b>2.30</b>	<b>0.014</b>	<b>0.981</b>
MSDESIS	1.71	9.71	1.39	2.91	0.0200	0.946
MCF-SMSIS	2.08	11.24	1.49	3.16	0.022	0.939
DUST3R (GT calib.)	71.82	100	231.10	231.64	4.08	0.00
DUST3R (DUST3R calib.)	81.63	100	231.10	231.64	4.08	0.00
RAFT	3.68	26.39	2.86	5.60	0.039	0.83
CATS-RAFTStereo	1.23	5.80	1.07	2.59	0.016	0.978
CATS-IGEV++	<u>1.21</u>	<u>5.52</u>	<u>1.01</u>	2.37	<u>0.015</u>	<u>0.980</u>

## 6 Conclusions

Camera augmentation bridges the gap between ideally conditioned synthetic datasets and the real surgical conditions where calibration is uncertain or unavailable. By introducing controlled geometric perturbations of camera parameters during training, it enables the adaptation of large-scale stereo matching models to the specific challenges and conditions of MIS. The proposed CATS demonstrates that robust and accurate stereo matching can be achieved even in the absence of calibration, while maintaining competitive performance when calibration is available. This capability stems from the in-depth analysis and modelling of stereo laparoscopes, allowing CATS to be rigorously tailored to the specific geometric variability encountered in surgical imaging. Beyond stereo matching, the concept of camera-aware image augmentation opens new perspectives for a broad range of computer vision tasks sensitive to camera geometry. By decoupling model performance from calibration constraints, this approach promotes the development of deep models that are more flexible, generalisable, and clinically deployable. As a direction for future work, we plan to test the CATS framework on a larger set of stereo surgical cameras, and to incorporate optical distortion models.

## References

- [1] Malhotra, S., Halabi, O., Dakua, S.P., Padhan, J., Paul, S., Palliyali, W.: Augmented reality in surgical navigation: a review of evaluation and validation metrics. *Applied Sciences* **13**(3), 1629 (2023)
- [2] Barcali, E., Iadanza, E., Manetti, L., Francia, P., Nardi, C., Bocchi, L.: Augmented reality in surgery: a scoping review. *Applied Sciences* **12**(14), 6890 (2022)
- [3] Xu, M., Guo, Z., Wang, A., Bai, L., Ren, H.: A review of 3D reconstruction techniques for deformable tissues in robotic surgery. In: *MICCAI*, pp. 157–167



528  
529 **Fig. 6** Qualitative model comparison in the uncalibrated scenario on samples from the RIS2017 and  
530 in-house datasets.

531 (2024)

- 533 [4] Scharstein, D., Szeliski, R.: A taxonomy and evaluation of dense two-frame stereo  
534 correspondence algorithms. *IJCV* **47**(1), 7–42 (2002)
- 535 [5] Lipson, L., Teed, Z., Deng, J.: RAFT-Stereo: Multilevel recurrent field transforms  
536 for stereo matching. In: *3DV*, pp. 218–227 (2021)
- 537 [6] Xu, G., Wang, X., Zhang, Z., Cheng, J., Liao, C., Yang, X.: IGEV++: Iterative  
538 multi-range geometry encoding volumes for stereo matching. *TPAMI* **47**(8), 7108–  
539 7122 (2025)
- 540 [7] Wen, B., Trepte, M., Aribido, J., Kautz, J., Gallo, O., Birchfield, S.: Foundation-  
541 Stereo: Zero-shot stereo matching. In: *CVPR*, pp. 5249–5260 (2025)
- 542 [8] Psychogyios, D., Mazomenos, E., Vasconcelos, F., Stoyanov, D.: MSDESIS: Mul-  
543 titask stereo disparity estimation and surgical instrument segmentation. *TMI*  
544 **41**(11), 3218–3230 (2022)
- 545 [9] Wu, R., He, C., Liang, P., Liu, Y., Huang, Y., Liu, W., Shu, B., Xu, P., Chang, Q.:  
546 MCF-SMSIS: multi-tasking with complementary functions for stereo matching  
547 and surgical instrument segmentation. *Computers in Biology and Medicine* **179**,  
548 549  
550  
551  
552

108923 (2024)	553
	554
[10] Li, L., Li, X., Ouyang, B., Mo, H., Ren, H., Yang, S.: Three-dimensional collision avoidance method for robot-assisted minimally invasive surgery. <i>Cyborg and Bionic Systems</i> <b>4</b> , 0042 (2023)	555
	556
	557
	558
[11] Narasimhan, S., Turkcan, M.K., Ballo, M., Choksi, S., Filicori, F., Kostic, Z.: Monocular 3D tooltip tracking in robotic surgery—building a multi-stage pipeline. <i>Electronics</i> <b>14</b> (10), 2075 (2025)	559
	560
	561
	562
[12] Preetha, C.J., Kloss, J., Wehrtmann, F.S., Sharan, L., Fan, C., Müller-Stich, B.P., Nickel, F., Engelhardt, S.: Towards augmented reality-based suturing in monocular laparoscopic training. In: <i>Medical Imaging</i> , p. 113150 (2020)	563
	564
	565
	566
[13] Banks, M.S., Read, J.C., Allison, R.S., Watt, S.J.: Stereoscopy and the human visual system. <i>Motion imaging journal</i> <b>121</b> (4), 24–43 (2012)	567
	568
	569
[14] Allan, M., Mcleod, J., Wang, C., Rosenthal, J.C., Hu, Z., Gard, N., Eisert, P., Fu, K.X., Zeffiro, T., Xia, W.: Stereo correspondence and reconstruction of endoscopic data challenge. <i>arXiv preprint arXiv:2101.01133</i> (2021)	570
	571
	572
[15] Hayoz, M., Hahne, C., Gallardo, M., Candinas, D., Kurmann, T., Allan, M., Sznitman, R.: Learning how to robustly estimate camera pose in endoscopic videos. <i>IJCARS</i> <b>18</b> (7), 1185–1192 (2023)	573
	574
	575
	576
[16] Allan, M., Shvets, A., Kurmann, T., Zhang, Z., Duggal, R., Su, Y.-H., Rieke, N., Laina, I., Kalavakonda, N., Bodenstedt, S.: 2017 robotic instrument segmentation challenge. <i>arXiv preprint arXiv:1902.06426</i> (2019)	577
	578
	579
	580
[17] Akui, N., Honma, S., Kanamori, I., Takahashi, S., Hanzawa, T., Fukaya, T., Karasawa, H., Kubota, T., Hashiguchi, T., Mochida, A.: Three-dimensional vision endoscope with position adjustment means for imaging device and visual field mask. <i>US Patent 5,577,991</i> (1996)	581
	582
	583
	584
	585
[18] Zhao, W., Mohr, C., DiMaio, S.: Stereo imaging system with automatic disparity adjustment for displaying close range objects. <i>US Patent 10,178,368</i> (2019)	586
	587
	588
[19] Breidenthal, R.S., Forkey, R.E., Smith, J., Volk, B.E.: Stereoscopic endoscope with virtual reality viewing. <i>US Patent 6,139,490</i> (2000)	589
	590
	591
[20] Mayer, N., Ilg, E., Hausser, P., Fischer, P., Cremers, D., Dosovitskiy, A., Brox, T.: A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In: <i>CVPR</i> , pp. 4040–4048 (2016)	592
	593
	594
[21] Teed, Z., Deng, J.: RAFT: Recurrent all-pairs field transforms for optical flow. In: <i>ECCV</i> , pp. 402–419 (2020)	595
	596
	597
[22] Wang, S., Leroy, V., Cabon, Y., Chidlovskii, B., Revaud, J.: DUST3R: Geometric	598

599 3d vision made easy. In: CVPR, pp. 20697–20709 (2024)  
600  
601 [23] Sharifian, R., Rabbani, N., Bartoli, A.: The RoDEM benchmark: evaluating the  
602 robustness of monocular single-shot depth estimation methods in minimally-  
603 invasive surgery. IJCARS **20**(6), 1215–1229 (2025)  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644