

**LETTER**

# Scale-preserving shape reconstruction from monocular endoscope image sequences by supervised depth learning

Takeshi Masuda<sup>1</sup>  | Ryusuke Sagawa<sup>1</sup> | Ryo Furukawa<sup>2</sup>  | Hiroshi Kawasaki<sup>3</sup>

<sup>1</sup>Artificial Intelligence Research Center, National Institute of Advanced Industrial Science and Technology (AIST), Tsukuba, Ibaraki, Japan

<sup>2</sup>Faculty of Engineering, Kindai University, Higashihiroshima, Hiroshima, Japan

<sup>3</sup>Faculty of Information Science and Electrical Engineering, Kyushu University, Fukuoka, Japan

**Correspondence**

Takeshi Masuda, Artificial Intelligence Research Center, National Institute of Advanced Industrial Science and Technology (AIST), Ibaraki, 305-8568, Japan.

Email: t.masuda@aist.go.jp

**Funding information**

Japan Society for the Promotion of Science (JSPS), KAKENHI, Grant/Award Numbers: JP18H04119, JP20H00611, JP21H01457; New Energy and Industrial Technology Development Organization, Grant/Award Number: JPNP20006

**Abstract**

Reconstructing 3D shapes from images are becoming popular, but such methods usually estimate relative depth maps with ambiguous scales. A method for reconstructing a scale-preserving 3D shape from monocular endoscope image sequences through training an absolute depth prediction network is proposed. First, a dataset of synchronized sequences of RGB images and depth maps is created using an endoscope simulator. Then, a supervised depth prediction network is trained that estimates a depth map from a RGB image minimizing the loss compared to the ground-truth depth map. The predicted depth map sequence is aligned to reconstruct a 3D shape. Finally, the proposed method is applied to a real endoscope image sequence.

## 1 | INTRODUCTION

Endoscopes are devices that allow non-invasive direct observation of internal structures and are used in the medical and industrial fields. In gastrointestinal examinations and surgeries, wired endoscopes are mainly used for the digestive tract, excluding the small intestine, while wireless capsule endoscopes are used for the examination of the entire digestive tract.

Most endoscopes are equipped with a monocular camera. Stereoscopic endoscopes are used to provide a 3D view to the surgeon in medical surgery, although they are not yet widely used due to operability limitations. Attempts are being made to reconstruct 3D information from stereoscopic endoscopes [1].

Reconstructing 3D structures from images has been an important topic in the history of computer vision. Recently, monocular 3D reconstruction has been realized using neural networks. Eigen et al. [2] developed a supervised monocular

depth estimation using CNNs, whose depth error measure is scale-invariant, so the estimated depth is relative. Monocular 3D reconstruction is possible in limited domains where training images are provided, such as street views in the KITTI dataset [3] and indoor scenes in the NYU dataset [4]. Most methods for obtaining 3D structure from monocular images provide relative depth estimation.

Methods such as structure from motion (SfM) and simultaneous localization and mapping (SLAM) are re-modelled by the learning framework [5] to reconstruct object shape and camera motion from image sequences. SfMLearner [6] modelled SLAM with two neural networks: DispNet for estimating disparity from a single image and PoseNet for estimating camera motion in a short-term image sequence. These two networks are jointly optimized by minimizing photometric consistency loss, which evaluates the change in pixel values before and after camera motion. Wang et al. [7] extended PoseNet with differentiable direct visual odometry (DVO), and Godard et al. [8] (monodepth2) integrated PoseNet with ResNet-18 and introduced a

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2023 The Authors. *Healthcare Technology Letters* published by John Wiley & Sons Ltd on behalf of The Institution of Engineering and Technology.

binary mask to remove irregular motion. These methods do not require the ground truth depth map for training and are referred to as unsupervised methods.

Even if the depth has scale ambiguity, it can be constrained to be consistent throughout the image sequence. The geometry consistency loss was introduced to penalize depth inconsistency between consecutive frames [9]. This work was further improved by pre-correcting for rotation [10] and by combining it with an additional relative monocular depth estimation module for high-density reconstruction that allows deformations [11].

SfMLearner [6] was adapted to endoscopic image sequences by EndoSfMLearner [12]. The geometry consistency loss was used to measure predicted depth consistency in addition to RGB value consistency evaluated as well by the photometric consistency loss. This study makes use of a capsule endoscope simulator called VR-Caps [13] to generate their own dataset for training and evaluation. The predicted depth maps still have scale ambiguity and they need to be rescaled appropriately for evaluation.

Predicting in the correct scale is beneficial for real medical applications. If a reference object is available, scale can be recovered by the post-processing, but it would be beneficial if absolute shape and position could be predicted online without any reference targets. DispNet [6] can be thought as an encoder-decoder model, and the encoder part can be replaced by various image encoder networks. Fang et al. [14] tested various encoders for supervised and unsupervised training of DispNet. Their conclusion is in short that supervised depth estimation is superior to unsupervised depth estimation.

This paper presents an absolute depth estimation method for monocular endoscopic image sequences. In general, monocular depth prediction is possible by assuming a specific scene domain. For our case, it is possible for endoscopic images of the digestive organs because the organs are usually in limited shape and colour, surely with individual differences, lighting and surface conditions. Compared to typical 3D reconstructed scenes such as KITTI and NYU, the difficulties in analyzing endoscopic scenes are featurelessness, wet specular reflections, uneven illumination, and severe turbulence and deformation. Therefore, using VR-Caps [13], we created a dataset that simulates a gastrointestinal endoscopy sequence with synchronized pairs of RGB images and ground truth depth maps. Using this dataset, we first train DispNet in a supervised manner using the ground truth depth maps with the correct scale, and then reconstruct the 3D shape by aligning and merging the predicted depth maps. Finally, the developed method was applied to real endoscopic image sequences.

The contributions of this paper are: first, we generated a synthetic dataset of the sequences of images with synchronized ground truth depth maps by an endoscope simulator, then we used it for training the DispNet to predict scale preserving depth maps from image sequences, and finally, we developed a SLAM algorithm to reconstruct the 3D shape from

the predicted depth maps, and applied it to real endoscope image sequences.

We first explain dataset preparation in Section 2, then the supervised training of DispNet in Section 3 and PoseNet in Section 4. In Section 5, a SLAM method to reconstruct the 3D shape is proposed, and we show the results of application to the real endoscope image sequences in Section 6.

## 2 | DATASET PREPARATION

For supervised training, the ground truth depth maps are needed in addition to the RGB images. Currently, it is still difficult to obtain reliable real depth maps synchronized with the image frames from endoscopes. Instead, we used VR-Caps [13], a Unity [15]-powered capsule endoscopy simulator, to generate computer-graphics datasets. Its graphical shape model is derived from a human CT scan and includes the major digestive organs from the stomach to the large intestine, with realistic textures applied. It features configurable camera, reflection, and illumination parameters, and can synchronously render a pair of realistic RGB images and depth maps by controlling camera position and pose.

In the original VR-Caps dataset [16], the image size is  $320 \times 320$  and the projection matrix and lens distortion are set to simulate a real capsule endoscope. In their dataset, we found that misalignments between the RGB images and the depth map is incomplete (Figure 1). We suppose that this misalignment is caused by capture timing delay due to online capturing, and we recreated our own dataset using the same system by off-line rendering with perfect synchronization. For generating our own dataset, the image size is kept the same as the original ( $320 \times 320$ ) for both the RGB images and depth maps, and the focal length is 159.45 pixels and the projection center is at the image center. We didn't apply lens distortion because it can be easily undistorted even for real endoscope images if calibrated.

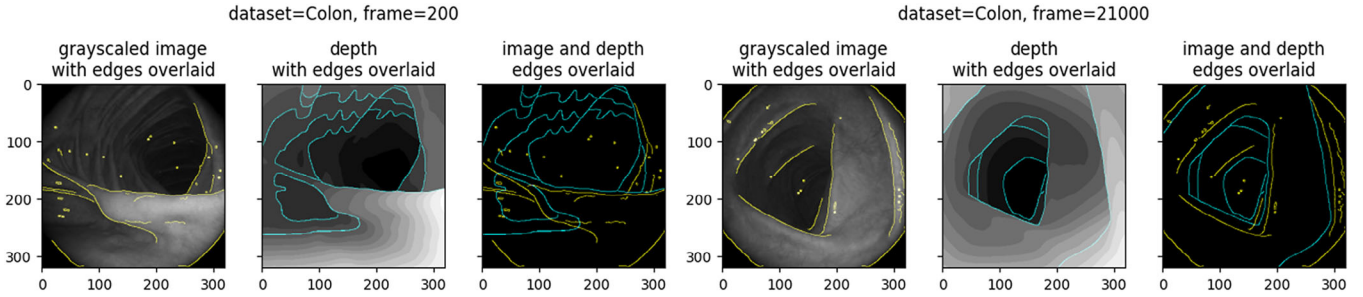
With these settings, we manually controlled the camera position and pose, and finally generated 17 sequences of 9714 frames in total, and each frame is composed of synchronized frames of a RGB-image and a depth map.

We split this dataset into two subsets: 13 sequences of 8249 frames for training and 4 sequences of 1465 frames for validation.

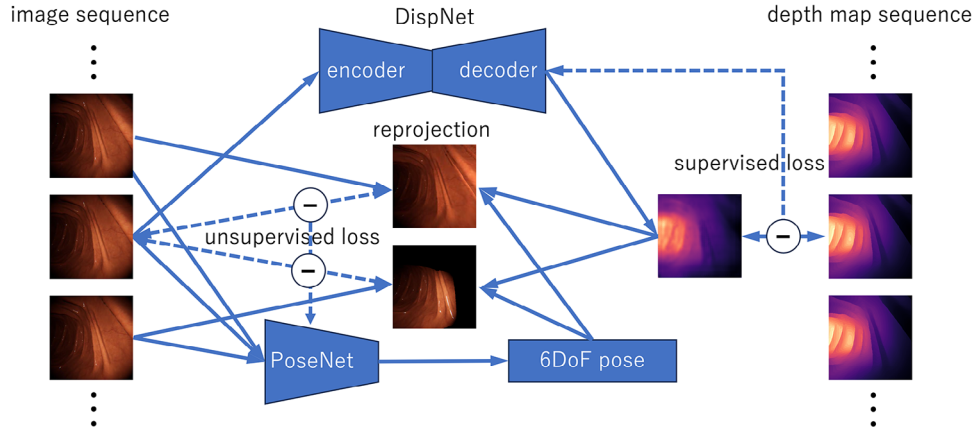
## 3 | SUPERVISED DEPTH ESTIMATION

### 3.1 | DispNet

DispNet is a neural network model that estimate a disparity map from a single RGB image (Figure 2). It consists of an encoder and a decoder, where the former extracts image features and the latter synthesizes the disparity map. The supervised DispNet [14] was trained on the ground truth depth maps and RGB images with validate various encoder structures. Following



**FIGURE 1** Two examples of misalignment of image and depth map in the original VR-Caps dataset [13]. In each example, left and middle show the RGB image and the depth map at the same frame in grey-scale. For clear visualization, we extracted the image edges by Canny operator, and overlaid them on the grey scaled image and the depth map in yellow and cyan respectively. These two-coloured edges are overlaid in right, which shows that the edges of image and depth map are not perfectly aligned. The amount and orientation of misalignment are not consistent over the whole sequence.



**FIGURE 2** Diagram of the method. We assume that we have synchronized sequences of RGB images and depth maps. It is composed of DispNet and PoseNet, and they are evaluated by the supervised and unsupervised losses.

many SfMLearner variants, the depth map  $D$  is obtained by taking the inverse of the disparity map  $d$  as  $D = 1/d$ .

The decoder part of DispNet consists of four-layers of up-convolution blocks. Each output of the up-convolution block is followed by a sigmoid function and a linear transformation, where  $d' = \alpha \text{sigmoid}(d) + \beta$  where  $d$  and  $d'$  are the disparity map before and after the transformation. The parameters  $\alpha$  and  $\beta$  are fixed to match the output range of the sigmoid function  $[-1, 1]$  to the target disparity range, which is commonly used by many successors of SfMLearner [6]. We replaced this to the softplus function  $d' = \text{softplus}(d) = \log(1 + e^d)$ . This ensures that the predicted disparity is always non-negative and eliminates the tuning of  $\alpha$  and  $\beta$  for the target dataset.

### 3.2 | Supervised training of DispNet

For obtaining the scale-preserving monocular depth estimator, we train DispNet by optimizing the loss function comparing the predicted depth map and the ground truth depth map [14]. We used the synthetic dataset of synchronized sequential pairs of RGB images and depth maps for training DispNet (Section 2).

As for the loss for the supervised depth estimation, we used the simple  $L_1$  loss:

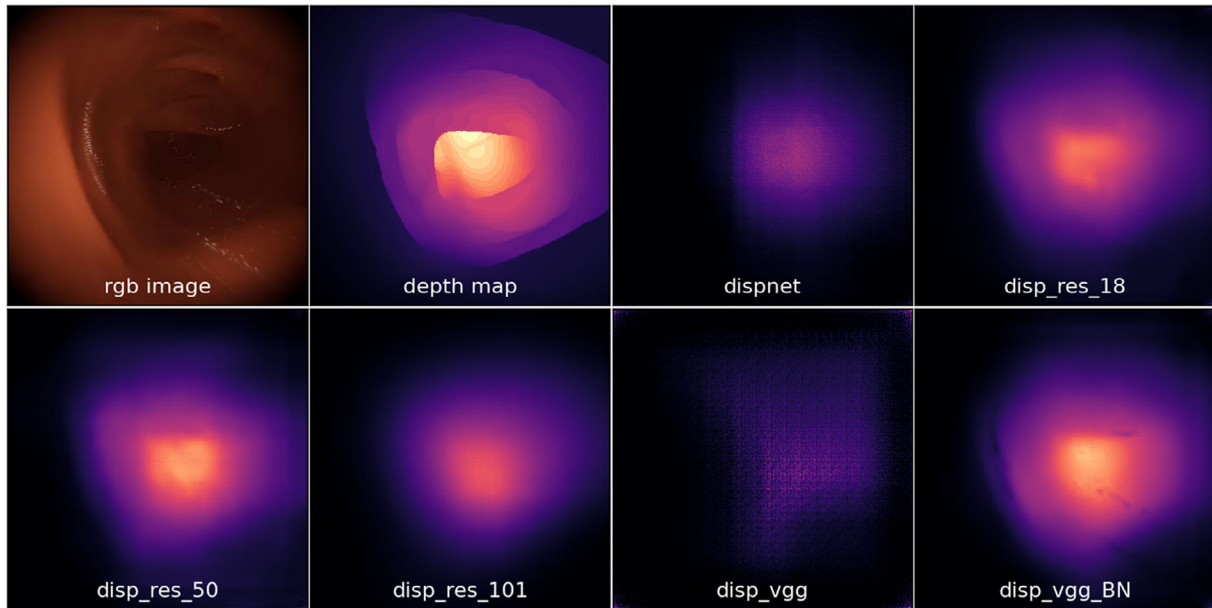
$$E_{\text{sup}} = \frac{1}{N} \sum_i |\hat{D}_i - \tilde{D}_i|, \quad (1)$$

where  $\tilde{D}_i$  and  $\hat{D}_i$  are the predicted and true depth values respectively at all corresponding pixels, and  $N$  is the total number of the pixels.

We trained DispNet from scratch with the batch size of 4 and the data augmentation of random horizontal flip and random 90-degree rotations. RGB values were normalized so that RGB channels have a mean of  $[0.5, 0.5, 0.5]$  and a standard deviation of  $[0.5, 0.5, 0.5]$ . The optimizer was Adam with the learning rate:  $10^{-7}$ , momentum: 0.9, beta: 0.999, no weight decay. We used a cloud-based computing system called ABCI [17].

### 3.3 | Ablation study of DispNet encoders

We compared the performance of various encoders of DispNet for depth prediction from endoscope images. Among the encoders tested in [14], we used the following five encoders



**FIGURE 3** The trained network was applied to the RGB images in the validation dataset to predict the depth map. From top left to bottom right, the input RGB image, the truth depth map, and the predicted depth maps of the various trained networks. All depth map colour maps are the same: black-magenta pixels are near and yellow-white pixels are far apart.

**TABLE 1** Results for the VR-Caps validation dataset. The columns represent, respectively, the network model, the number of model parameters to train, the minimum training mean absolute error (MAE) and the number of epochs in which it occurred, the validation MAE (lower is better) and correlation coefficient (CC: higher is better), and training time. The best model is shown in bold and the second best model in italics. The learning epochs for DispNet was 500; the other models were 1000. The units for MAE are meters.

Network	#params	Least MAE	MAE	CC	time
dispnet	31,596,900	1.208e-2@111	0.01172	0.328	10:14
disp_res_18	14,330,244	5.955e-3@907	0.00568	0.872	20:29
disp_res_50	32,523,140	5.724e-3@992	0.00540	0.873	35:16
disp_res_101	62,907,268	7.584e-3@919	0.00728	0.782	57:58
disp_vgg	143,507,564	1.158e-2@30	0.01120	0.297	39:27
disp_vgg_BN	143,516,012	5.665e-3@919	<b>0.00539</b>	<b>0.883</b>	42:21

for DispNet: `dispnet`, `disp_res_18`, `disp_res_50`, `disp_res_101`, `disp_vgg`, `disp_vgg_BN`, which stand for the original DispNet [2, 6], DispNet with ResNet-18/50/101 and VGG without/with batch normalization, respectively.

The optimization results are summarized in Table 1. The validation loss of the `dispnet` and `disp_vgg` were minimized in an early stage of the optimization process, which may mean that they are suffered from the over fitting. In terms of the MAE and the correlation coefficient, `disp_vgg_BN` was the best, but `disp_res_50` and `disp_res_18` followed with very minor differences.

The predicted depth maps are shown in Figure 3. We can understand that the predictions of `disp_res_18`, `disp_res_50`, and `disp_vgg_BN` are visually close to

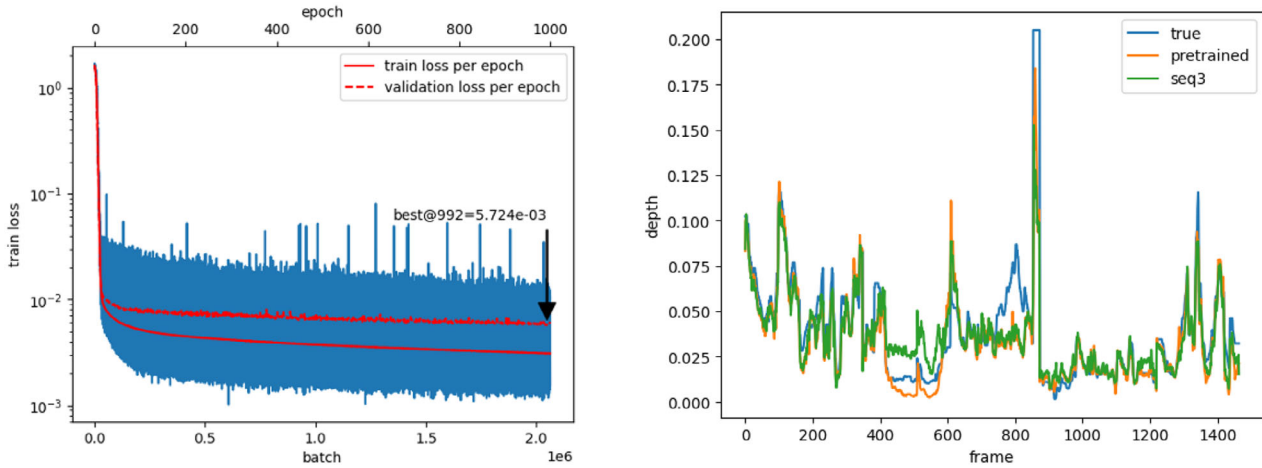
the true depth map, which supports the previous insights of error metrics. When we carefully inspect the predicted depth map, there are some noisy dots in `disp_vgg_BN` cases that don't exist in the ground truth depth map. Although it does not appear in Table 1, from a viewpoint of stability, we consider that ResNet based encoders (`disp_res_18/_50`) are more stable than VGG based encoders.

Figure 4 shows the supervised training procedure and the result of DispNet with ResNet-50 encoder (`disp_res_50`). Comparison of the true and predicted depth values shows that the prediction of absolute depth values was successful.

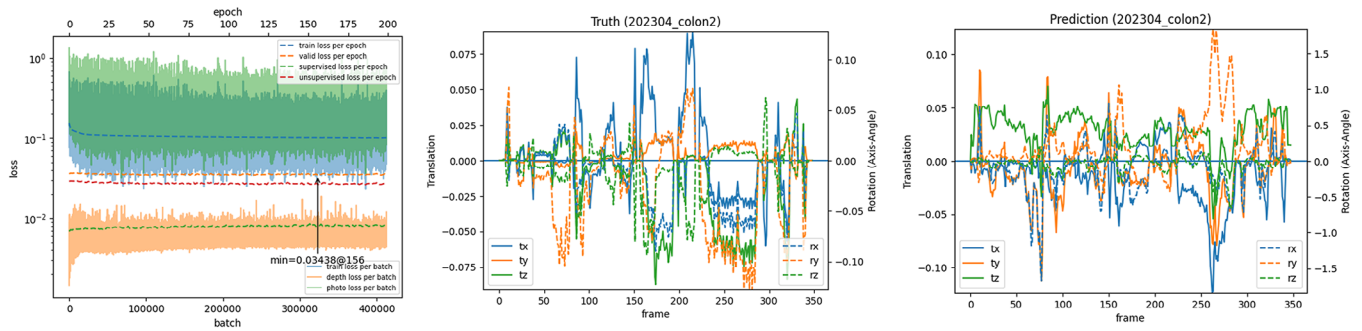
## 4 | POSE ESTIMATION

### 4.1 | PoseNet

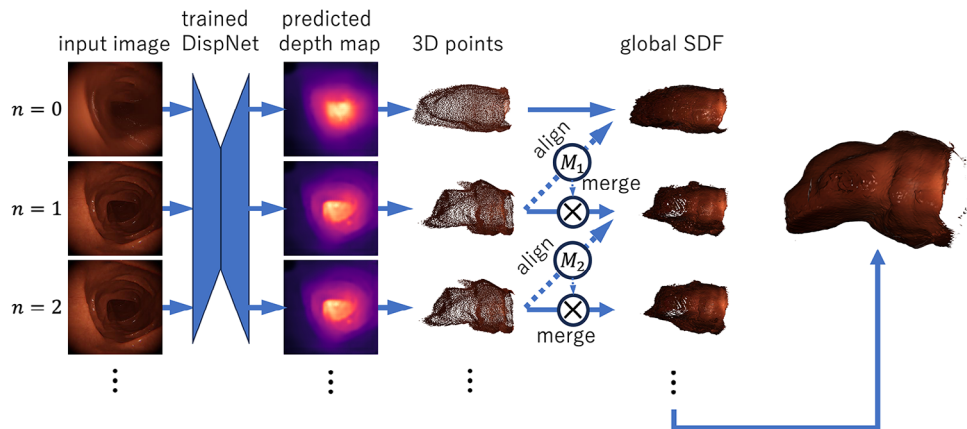
PoseNet is a neural network that predicts a 6-DoF 3D relative rigid motion parameters from a sequence of input images frames [18], and there are many varieties of implementations. For the target image, we take its preceding and succeeding images as the reference images to form a sequence. We used two types of PoseNet with different sequence length: one uses the target image and all the reference images as the input, and the other uses the target image and only one of the reference images, which are respectively referred by ‘‘PoseExpNet’’ and ‘‘PoseResNet’’ as appeared in the source codes of [6, 14] and [8, 12]. The former is implemented as a pure CNN with 1.6M parameters, and the latter uses ResNet18 as the feature detector followed by a CNN with 13M parameters in total.



**FIGURE 4** Left: plot of training and validation losses of DispNet with ResNet-50 encoder (`disp_res_50`). Right: The values of the true and predicted depth map at the image center (160,160) through the sequence.



**FIGURE 5** Training results of PoseNet. Left: plot of training and validation losses during training. Middle and right: the relative motion between successive frames.



**FIGURE 6** Graphical representation of 3D shape reconstruction (Algorithm 1). With the trained DispNet, an absolute depth map is predicted from each input image. The depth maps are further aligned and merged to the global shape.

There are also varieties of 6-DoF  $SE(3)$  pose parameterization used for SfMLearner descendants: a translation vector and a quaternion for rotation [6], a translation vector and Euler angles [12], a translation vector and a rotation vector (axis-angle representation or Rodrigues' formula) [7, 8, 10].

In the followings, we show the result of PoseExpNet that takes the next and previous images from the target image as the reference images, where the sequence length is 3, with the pose representation of the exponential map [19].



**FIGURE 7** 3D shape reconstruction results. The reconstructed shape from 50 frames of the true depth map (left) and the predicted depth map (middle), and the true poses and the estimated poses are compared (right). The rotation parameters are in the axis-angle representation. Back faces are culled off for visualizing meshes.

#### ALGORITHM 1 Pseudocode for our SLAM algorithm

---

**Require:**  $\{D_n\}$  ( $0 \leq n \leq N - 1$ ),  $K \triangleright$  A sequence of depth maps and the intrinsic matrix

- 1:  $M_0 \leftarrow I_{4 \times 4} \triangleright$  the pose of the first frame is fixed to the identity matrix
- 2:  $V \leftarrow \emptyset \triangleright$  set the global SDF voxels
- 3: **for**  $n \leftarrow 0, N - 1$  **do**
- 4:  $P_n \leftarrow \text{convertToPointSet}(D_n, K) \triangleright$  convert a depth map to a point set
- 5: **if**  $n > 0$  **then**
- 6:  $Q \leftarrow \text{extractPointSet}(V) \triangleright$  extract a point set from the global SDF voxels
- 7:  $M_n \leftarrow \text{ICP}(P_n, Q, M_{n-1}) \triangleright$  optimize the pose  $M_n$  from  $M_{n-1}$  such that  $Q \approx M_n(P_n)$
- 8: **end if**
- 9:  $V \leftarrow V \cup M_n(P_n) \triangleright$  merge the transformed point set to the global SDF voxels
- if**  $n > 0$  and  $n \bmod N_1 = 0$  **then**  $\triangleright$  global alignment at every  $N_1 = 50$  iterations
- 10:  $Q \leftarrow \text{extractPointSet}(V) \triangleright$  extract a point set from the global SDF voxels
- 11: **for**  $m \leftarrow 1, n$  **do**
- 12:  $M_m \leftarrow \text{ICP}(P_m, Q, M_m) \triangleright$  update pose  $M_m$
- 13: **end for**
- 14:  $V \leftarrow \emptyset \triangleright$  reset the global SDF voxels
- 15: **for**  $m \leftarrow 0, n$  **do**
- 16:  $V \leftarrow V \cup M_m(P_m) \triangleright$  update global SDF voxels
- 17: **end for**
- 18: **end if**
- 19: **end for**
- 20: **end for**

---

## 4.2 | Unsupervised training of PoseNet

With the supervised DispNet, we train PoseNet in an unsupervised manner using its prediction. We used the loss function that is a weighted sum of the supervised and unsupervised losses:

$$E_{\text{total}} = w_{\text{sup}} E_{\text{sup}} + w_{\text{unsup}} E_{\text{unsup}},$$

where  $E_{\text{sup}}$  and  $E_{\text{unsup}}$  signify respectively the supervised depth loss (Equation (1)) and the unsupervised loss. The unsupervised loss is the  $L_1$  (the mean absolute error) of RGB values between the reprojected reference image and the target image, which is identical to the photometric consistency loss of [6, 14]. The weights  $w_{\text{sup}}$  and  $w_{\text{unsup}}$  are for balancing these losses, and we set  $w_{\text{sup}} = 1$  and  $w_{\text{unsup}} = 1$ .

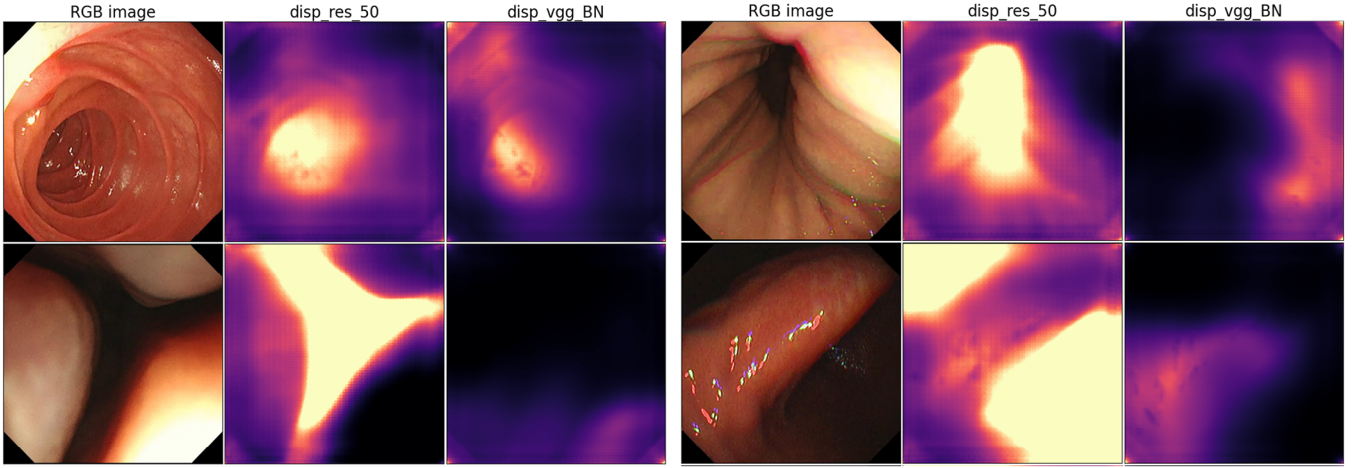
We used DispNet with ResNet-50 encoders whose parameters were pre-trained in Section 3.2. We set the learning rates for DispNet and PoseNet to  $10^{-7}$  and  $10^{-4}$  respectively, and the optimizer settings are the same as Section 3.2. The optimization process was iterated for 200 epochs.

Figure 5 shows that the pose prediction is still imperfect and much more optimization is needed. The possible causes are the dataset, where the endoscope images are generally uniform and featureless, and our choice of PoseNet, pose parameterization and loss function.

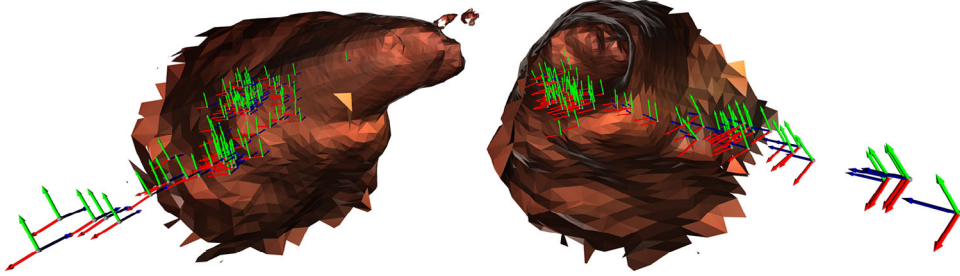
## 5 | 3D SHAPE RECONSTRUCTION

With the trained PoseNet, we can predict the relative pose between the image frames, but the pose prediction is not reliable. In the original SfMLearner [6], the authors reported their pose estimation by the PoseNet outperformed the ORB-SLAM [9], but in many its variants, the output of PoseNet is refined by the ICP algorithm [12] or by the other SLAM methods like ORB-SLAM [9].

In our case, PoseNet did not predict the pose with sufficient quality, and we use the ICP algorithm [20] for more reliable shape reconstruction. Instead applying the ICP algorithm to



**FIGURE 8** The example cases where evident difference can be observed in comparing the prediction of the DispNet with two different encoders: `disp_res_50` and `disp_vgg_BN`. The depth colourmap is common to all depth maps: black-magenta pixels are near and yellow-white pixels are far.



**FIGURE 9** The reconstructed 3D shape and the camera poses rendered from two viewpoints. Each triplet of arrows represents the estimated pose of the camera. Back faces are culled off for visualizing meshes.

successive frames, we use the global signed distance field (SDF) voxels as the global (canonical) shape, and align each frame to it [21, 22]. We developed a SLAM algorithm (Algorithm 1, Figure 6) for reconstructing the global shape and the camera poses implemented by the Open3D Python packages [23].

The function “convertToPointSet” converts a depth map to a point set, where a depth value  $D$  at  $(u, v)$  of a depth map is converted to a 3D point by  $\mathbf{p} = D \cdot \mathbf{K}^{-1} \cdot (u, v, 1)^T$ . The transformation  $M(P)$  represents a 3D rigid motion whose rotation and translation are respectively  $R$  and  $\mathbf{t}$ , where a point  $\mathbf{p} \in P$  is transformed to a point  $\mathbf{p}'$  by  $\mathbf{p}' = R\mathbf{p} + \mathbf{t}$ . The SDF voxels are implemented by `open3d.integration.ScalableTSDFVolume`, and its `integrate` and `extract_point_cloud` functions are used for “extractPointSet” and merging of point sets at lines 9 and 17 respectively in Algorithm 1.

Figure 7 shows the reconstruction result of this algorithm from the depth sequence of 50 frames predicted from the synthetic images. The predicted depth value is used only in the central region where pixels are located within the circle circumscribing the image boundary. The plots of poses estimated from the true depth is almost identical to the ground truth pose, and the 3D shape reconstructed from the predicted depth is similar to that from the true depth, which shows that the proposed algorithm worked correctly. For pose estimation from

the predicted depth map, the rotation estimation is not accurate because the object shape is nearly cylindrical and no colour information was used.

## 6 | APPLICATION TO REAL DATA

### 6.1 | Depth prediction

We applied the trained supervised DispNet to a real dataset that contains 13 RGB image sequences of digestive organs. There is no ground truth depth map that can be used for evaluation in this dataset. We estimated the camera parameters including the lens distortion from small portion of the sequence where the COLMAP [24] algorithm could work on, and rectify the images to adapt our DispNet. We chose two encoders: `disp_res_50` and `disp_vgg_BN` for comparison, and we found that there are frames with significant difference between these two encoders: selected cases are shown in Figure 8. The performances of these two encoders were similar as observed in Table 1 in terms of the error metrics, but considering the adaptation capability to the real data, `disp_res_50` works more stably than `disp_vgg_BN`. Due to lack of real endoscope images with ground truth depth map, we could not apply additional training for adaptation.

## 6.2 | 3D Shape reconstruction

We applied our 3D reconstruction algorithm (Algorithm 1) on the predicted depth maps. Figure 9 shows the result of reconstruction from 110 images. Due to lack of the ground truth depth map, we cannot evaluate this result like as in Section 6.1, but this figure shows that a reasonable cylindrical 3D shape was reconstructed by the proposed algorithm.

## 7 | CONCLUSIONS

This paper presents supervised learning of depth prediction networks from monocular RGB images with absolute depth preservation. Synchronous sequences of RGB images and depth maps were generated using an endoscope simulator. The performance of various DispNet models was compared, and as far as we tested, ResNet-based encoder performed the best. By training DispNet with the ground truth depth maps, the absolute depth maps were predicted from RGB images. According to interviews with clinicians, reconstructed shape accuracy should be at least 5 mm or less for practical use, and our results largely satisfied this requirement. We also developed a SLAM algorithm based on the ICP algorithm and SDF to align and integrate the predicted absolute depth maps to form a scale-preserving 3D shape model. These methods were applied to actual endoscopic image sequences.

We hope to improve our study to perform better in a variety of endoscopic scenes. Since it is difficult to obtain datasets of real endoscope RGB images with synchronized ground truth depth maps, we trained with synthetic datasets, but we need datasets with much varieties. For PoseNet, the performance for endoscopic images needs to be verified, since results of sufficient quality were not obtained. The proposed SLAM algorithm worked properly, but to improve its performance, it is necessary to use colour and texture information in addition. Also, we would like to extend the SLAM algorithm to handle more realistic cases including deformations and turbulence.

## NOMENCLATURE

SfM	Structure from motion
SLAM	Simultaneous localization and mapping
CNN	Convolutional neural network
MAE	Mean absolute error
ICP	Iterative closest point

## AUTHOR CONTRIBUTIONS

**Takeshi Masuda:** Investigation; methodology; software; visualization; writing—original draft; writing—review and editing. **Ryusuke Sagawa:** Data curation; funding acquisition; project administration; resources; supervision. **Ryo Furukawa:** Data curation; funding acquisition; project administration; supervision. **Hiroshi Kawasaki:** Funding acquisition; project administration; supervision.

## ACKNOWLEDGMENTS

This work was supported by JSPS KAKENHI Grant Numbers JP20H00611, JP18H04119, JP21H01457. This paper is based on results obtained from a project, JPNP20006, subsidized by the New Energy and Industrial Technology Development Organization (NEDO).

## CONFLICT OF INTEREST STATEMENT

The authors declare no conflicts of interest.

## DATA AVAILABILITY STATEMENT

Data that support the findings of this study are partially available from the corresponding author upon reasonable request. Open repository is not prepared yet.

## ORCID

*Takeshi Masuda*  <https://orcid.org/0000-0003-3449-3424>

*Ryo Furukawa*  <https://orcid.org/0000-0002-2063-1008>

## REFERENCES

- Wang, Y., Long, Y., Fan, S.H., Dou, Q.: Neural rendering for stereo 3D reconstruction of deformable tissues in robotic surgery. In: 25th International Conference on Medical Image Computing and Computer Assisted Intervention (MICCAI), pp. 431–441. Springer, Cham (2022)
- Eigen, D., Puhrsch, C., Fergus, R.: Depth map prediction from a single image using a multi-scale deep network. In: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N., Weinberger, K.Q. (eds.) Advances in Neural Information Processing Systems, vol. 27, Curran Associates, Red Hook, NY (2014)
- Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the KITTI vision benchmark suite. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition, pp. 3354–3361. IEEE, Piscataway, NJ (2012)
- Silberman, N., Hoiem, D., Kohli, P., Fergus, R.: Indoor segmentation and support inference from RGBD images. In: Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., Schmid, C. (eds.) Computer Vision – ECCV 2012, pp. 746–760. Springer, Berlin, Heidelberg (2012)
- Tateno, K., Tombari, F., Laina, I., Navab, N.: CNN-SLAM: Real-time dense monocular SLAM with learned depth prediction. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6565–6574. IEEE, Piscataway, NJ (2017)
- Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6612–6619. IEEE, Piscataway, NJ (2017)
- Wang, C., Miguel Buenaposada, J., Zhu, R., Lucey, S.: Learning depth from monocular videos using direct methods. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2022–2030. IEEE, Piscataway, NJ (2018)
- Godard, C., Mac Aodha, O., Firman, M., Brostow, G.J.: Digging into self-supervised monocular depth prediction. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 3827–3837. IEEE, Piscataway, NJ (2019)
- Bian, J.W., Zhan, H., Wang, N., Li, Z., Zhang, L., Shen, C., et al.: Unsupervised scale-consistent depth learning from video. *Int. J. Comput. Vis.* 129, 2548–2564 (2021). doi: <https://doi.org/10.1007/s11263-021-01484-6>
- Bian, J.W., Zhan, H., Wang, N., Chin, T.J., Shen, C., Reid, I.: Auto-rectify network for unsupervised indoor depth estimation. *IEEE Trans. Pattern Anal. Mach. Intell.* 44(12), 9802–9813 (2021). doi: <https://doi.org/10.1109/TPAMI.2021.3136220>
- Sun, L., Bian, J.W., Zhan, H., Yin, W., Reid, I., Shen, C.: SC-DepthV3: robust self-supervised monocular depth estimation for dynamic scenes. *arXiv:2211.03660* (2022)



12. Ozyoruk, K.B., Gokceler, G.I., Bobrow, T.L., Coskun, G., Incetan, K., Almalioglu, Y., et al.: EndoSLAM dataset and an unsupervised monocular visual odometry and depth estimation approach for endoscopic videos. *Med. Image Anal.* 71, 102058 (2021). doi: <https://www.sciencedirect.com/science/article/pii/S1361841521001043>
13. İncetana, K., Celikb, I.O., Obeida, A., Gokcelera, G.I., Ozyoruka, K.B., Almalioglu, Y., et al.: VR-Caps: a virtual environment for capsule endoscopy. *Med. Image Anal.* 70, 101990 (2021)
14. Fang, Z., Chen, X., Chen, Y., Van Gool, L.: Towards good practice for CNN-based monocular depth estimation. In: *Proceedings of 2020 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1091–1100. IEEE, Piscataway, NJ (2020)
15. Haas, J.K.: A history of the Unity game engine, Dissertation, Worcester Polytechnic Institute (2014)
16. VR-Caps: a virtual environment for capsule endoscopy. <https://github.com/CapsuleEndoscope/VirtualCapsuleEndoscopy.git> (2020). Accessed 17 July 2023
17. Ai bridging cloud infrastructure (abci). <https://docs.abci.ai/en/>. Accessed 17 July 2023
18. Kendall, A., Grimes, M., Cipolla, R.: PoseNet: A convolutional network for real-time 6-DOF camera relocalization. In: *IEEE International Conference on Computer Vision (ICCV)*, pp. 2938–2946. IEEE, Piscataway, NJ (2015)
19. Murray, R.M., Li, Z., Sastry, S.S.: *A Mathematical Introduction to Robotic Manipulation*. CRC Press, Boca Raton, FL (1994)
20. Besl, P.J., McKay, N.D.: A method for registration of 3-D shapes. *IEEE Trans. Pattern Anal. Mach. Intell.* 14(2), 239–256 (1992)
21. Masuda, T.: Registration and integration of multiple range images by matching signed distance fields for object shape modeling. *Comput. Vis. Image Understanding* 87(1–3), 51–65 (2002)
22. Newcombe, R.A., Izadi, S., Hilliges, O., Molyneaux, D., Kim, D., Davison, A.J., et al.: KinectFusion: real-time dense surface mapping and tracking. In: *2011 10th IEEE International Symposium on Mixed and Augmented Reality*, pp. 127–136. IEEE, Piscataway, NJ (2011)
23. Zhou, Q.Y., Park, J., Koltun, V.: Open3D: a modern library for 3D data processing. arXiv:1801.09847 (2018)
24. Schönberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4104–4113. IEEE, Piscataway, NJ (2016)

**How to cite this article:** Masuda, T., Sagawa, R., Furukawa, R., Kawasaki, H.: Scale-preserving shape reconstruction from monocular endoscope image sequences by supervised depth learning. *Healthc. Technol. Lett.* 11, 76–84 (2024). <https://doi.org/10.1049/htl2.12064>