

Mobile-based 3D modeling: an in-depth evaluation for the application in indoor scenarios

Martin De Pellegrini ^{1,†,*}, Lorenzo Orlandi ^{1,‡}, Daniele Sevegnani ¹ and Nicola Conci ^{2,*}

¹ ARCODA s.r.l., (Trento, Italy)

² University of Trento, Italy

* Correspondence: martin.depellegrini@arcoda.it (M.DP.), nicola.conci@unitn.it (NC)

‡ These authors contributed equally to this work.

Abstract: Indoor environment modeling has become a relevant topic in several applications fields including Augmented, Virtual, and Extended Reality. With the digital transformation, many industries have investigated the possibility to generate detailed models of an indoor environment allowing the viewers to navigate through it, and mapping surfaces so as to insert virtual elements overlaid to the real scene. The scope of the paper is twofold. We first review the existing state-of-the-art (SoA) of learning-based methods for 3D scene reconstruction based on Structure From Motion (SFM) that predict depth maps and camera poses from a video stream. We then present an extensive evaluation using a recent SoA network, with particular attention to the capability of generalizing on new unseen data of indoor environments. The evaluation was conducted based using as a baseline metric the Absolute relative (AbsRel) measure on the depth map prediction.

Keywords: Computer Vision; 3D Reconstruction; Deep Learning; Indoor; Digital Twin; Point Cloud.

1. Introduction

The ability of sensing the 3D space using single cameras is a widely investigated topic in image processing and computer vision. Several solutions have been developed over the years to ensure a reliable reconstruction of the observed environment, adopting both traditional image processing [1][2][3], as well as more up-to-date learning approaches [9][34]. In fact, 3D sensing and reconstruction is a necessary building block behind a number of technologies in industry, including robotics, landslide mapping, gaming, mixed reality, archaeology, medicine, to name a few [4][5][6]. Despite the many efforts spent by the research community in providing progressively more accurate models capable and sensing and reconstructing a 3D environment, a number of challenges remains still open. In fact, the acquisition of 3D information can serve multiple purposes, and can be used in real-time in a multi-sensorial context, as seen in robots or, more in general, autonomous systems. This often implies that the visual information is only one among the multiple inputs to a localization and navigation system. In such conditions, the potential errors emerging from inaccuracies and/or wrong reconstruction of portions of the environment are often compensated and mitigated thanks to the presence of additional sensing devices. Vice versa, in a more restrictive context, in which multi-modal equipment is not a viable option, 3D reconstruction is performed using the visual information solely, thus requiring high resolution images for better feature detection, and accurate

Citation: Lastname, F.; Lastname, F.; Lastname, F. Title. *Journal Not Specified* **2021**, *1*, 0. <https://doi.org/>

Received:

Accepted:

Published:

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Copyright: © 2022 by the authors. Submitted to *Journal Not Specified* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

34 camera calibration with distortion correction in order to generate the 3D model,
35 consisting of a sparse or dense point cloud.

36 In this paper, we present an in-depth evaluation of a robust state-of-the-
37 art method for depth estimation, which is used as the core element for 3D
38 reconstruction applications. In particular, we focus our research on the indoor
39 scenario, in which we expect the user to collect data using an arbitrary camera,
40 and following subjective criteria. In other words, the acquisition is not conducted
41 following a rigorous path in scanning the environment, thus not imposing any
42 constraint on the user's side. Such conditions are indeed very common, and
43 cover a wide spectrum of application scenarios, often involving on-the-field
44 workers, which rely on such augmented/extended reality tools for inspection
45 and maintenance operations.

46 The paper is structured as follows: in Section 2 we present some recent
47 relevant related work; Section 3 discusses the motivation of this work and the
48 main contributions; in Section 4 we focus on the validation pipeline we have
49 envisaged, describing the methodology and the metrics used. In Section 5 the
50 achieved results are presented and discussed. Final remarks and conclusions are
51 drawn in Section 6.

52 2. Related Work

53 In the following paragraphs, we report the most relevant works presented
54 in the SoA, starting from the traditional Structure from Motion algorithm and
55 surveying the most recent developments based on deep-learning. Structure
56 from Motion (SfM) [28] allows the estimation of the three-dimensional structure
57 of objects and environments based on the motion parallax that describes the
58 appearance changes of an object when the observer's viewpoint changes. By
59 doing so, it is possible to infer the 3D structure of a target, and retrieve the
60 distance from the camera to generate a 3D representation. Another basic principle
61 of SfM is the stereo photogrammetry triangulation used to calculate the relative
62 position of points from stereo pairs. More in general, SfM is required to solve
63 three main problems. Firstly (i) it is required to find correspondences between
64 the images and measure the distances between the features extracted with respect
65 to the two image planes. Typically, SIFT [32] features are used in this phase due
66 to their robustness against change in scale, large variation of view point and
67 challenging conditions such as different illumination and partial occlusions; as a
68 second step, (ii) the camera position associated to each of the images processed
69 is computed, via bundle adjustment (BA), to calculate and optimize 3D structure,
70 camera pose and intrinsic calibration; lastly, (iii) generate a 3D dense point cloud
71 by using the camera parameters to back project the points computed before on
72 the 3D space, also called *multi view stereo matching*.

73 Traditional 3D reconstruction algorithms require to perform heavy operations
74 and despite the proven effectiveness of these methods, they rely on high
75 quality images as input. This may introduce some or limitations when it comes
76 to process complex geometries, occlusions and low-texture areas. Such issues
77 have been partially tackled replacing traditional feature and geometry-based
78 approaches with deep learning. In particular, some stages of the traditional 3D
79 reconstruction pipeline have been rethought following a deep learning-based for-
80 mulation. To this aim, we present some of the methods explored for the purpose

81 of our research, which implement the principles of SfM using Convolutional
 82 Neural Networks (CNNs). One of the most relevant works exploiting neural net-
 83 works for depth estimation is DispNet [7]. DispNet is used for single view depth
 84 prediction. It is composed by an initial contracting stage, made of convolutional
 85 layers, followed by up sampling to perform deconvolutions, convolutions and
 86 computation of the loss function. Features from the contracting part are sent to
 87 the corresponding layer in the expanding portion. The network operates with a
 88 traditional encoder-decoder architecture with skip connections and multi-scale
 89 side prediction. The DispNet architecture is reported for convenience in Figure 1.

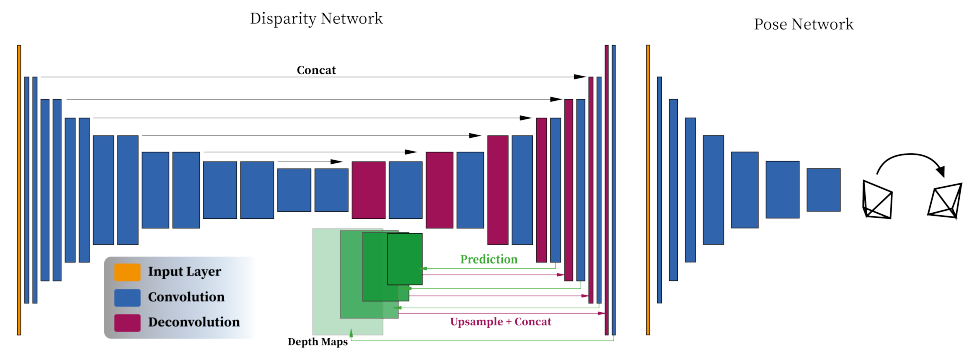


Figure 1. Illustration of the architecture of the Disparity estimation Network (DispNet) with encoder-decoder layout, and Pose estimation Network. Additional details in terms of the size of each layer can be found in the original manuscript.

90 Many solutions have been developed employing Convolutional Neural
 91 Networks (CNNs) for the task of estimating the depth information. Some of
 92 them are used for stereo view synthesis such as DeepStereo [8], which learns
 93 how to generate new views from single images in order to recreate a synthetic
 94 stereoscopic system where the underlying geometry is represented by quantized
 95 depth plane. Similarly, Deep3D [9] implements CNNs to convert 2D video into
 96 3D sequences such as Anaglyph for 3D movie or Side-by-Side view for Virtual
 97 Reality (VR) applications. In this case the scene geometry is represented by
 98 probabilistic disparity maps. As well as Deep3D, other methods are following the
 99 recent research in learning three-dimensional structure from single view. Some
 100 of them introduced supervision signals such as in the work proposed by Garg
 101 *et al.* [33]. The authors propose a supervision consisting of a calibrated stereo
 102 twin for single-view depth estimation. The recent trends in depth estimation
 103 aim for unsupervised or self-supervised learning from video sequences. These
 104 methods work well in the task of inferring the scene geometry and ego-motion
 105 (similarly to SfM), but in addition they show great potential for other tasks such
 106 as segmentation, object motion mask prediction, tracking and other levels of
 107 semantics (please refer to [18][19][20][21][22][23]).

108 Among the unsupervised/self-supervised methods, three important re-
 109 searches have been conducted by Vijayanarasimhan *et al.* [34], Zhou *et al.* [10]
 110 and Bian *et al.* [11]. These approaches implement two sub networks: the first one
 111 focuses on single view depth prediction, and the second one is used for camera
 112 pose estimation in support to the depth network, so as to replicate a pseudo
 113 stereo vision setting. These implementations mostly differ on the loss function,
 114 which is applied as supervision signal. In terms of performances the methods
 115 achieve state-of-the-art scores on the KITTI [31] and Cityscapes [35] datasets.

Ref.	Method	Indoor	Dataset	Note
[34]	SfM Net	✗	KITTI[31] & Cityscapes[35]	O
[10]	SfM Learner	✗	KITTI[31] & Cityscapes[35]	O
[11]	SC-SfM Learner	✗	KITTI[31] & Cityscapes[35]	O
[12]	Indoor SC-SfM Learner	✓	NYUv2[27]	R

Table 1: Methods from literature for depth estimation from video sequences. In the column **Note** symbols (O) and (R) refer to Original and Rectified Training data.

116 3. Motivation and Contribution

117 Despite the proven effectiveness in street mapping contexts, the previous
 118 methods do not perform well when it comes to infer the 3D structure of indoor
 119 environments, and also by training the network with indoor RGB-D datasets, it
 120 does not allow to achieve satisfactory results, as also mentioned in [12]. Indeed,
 121 DispNet aims to learn the disparity between frames and due to the nature of hand-
 122 recorded sequences, typical of indoor data collection, the spatial relationship
 123 between adjacent frames might be of pure rotation, leading to a disparity equal
 124 to zero. More in detail, it has been demonstrated that the estimation of the depth
 125 map is strictly related to a dominance of translation with respect to rotations
 126 in the video sequences acquisition. In fact, previous implementations have
 127 been tested on datasets like KITTI [31], where the camera configuration and the
 128 forward motion did not give evidence to this issue. A research conducted by
 129 Bian *et al.* [12] has proven the existence of such limitation of the DispNet and
 130 proposes a weak rectification algorithm to pre-process indoor datasets before
 131 training the network. The authors have applied the rectification on the NYUv2
 132 [27][25] dataset used to train the network and tested the generalization capability
 133 on the *7Scene* dataset [17]. Since the generalization was evaluated on one dataset
 134 only, we aim to provide additional benchmarks evaluating other RGB-D datasets
 135 and comment on the network generalization capability.

136 In summary the main contributions of the paper are:

- 137 • We provide additional benchmarks for the network proposed by Bian *et al.*
 138 *in order to allow a better understanding on the network generalization*
 139 *performances.*
- 140 • We analyze the network generalization capability in connection with the
 141 *statistics of the scene, from which the depth has to be estimated. We com-*
 142 *pute the depth standard deviation from depth ground truth to describe the*
 143 *amount of depth information that the network has to estimate, and then*
 144 *discuss how the generalization is related to this parameter.*

145 4. Materials and Methods

146 As anticipated in the previous section, the results and evaluation that are
 147 presented in the following paragraphs are based on the work by Bian *et al.* [11]
 148 [12]. Here, the network model is pre-trained on ImageNet [13] using ResNet-18
 149 [14] in substitution to the depth and pose sub networks. Next, a fine-tuning
 150 on the rectified NYUv2 (Figure 3) [27][25] dataset is applied. Differently from
 151 the other architectures, the framework has been developed to overcome the
 152 scale ambiguity in [10], but it preserves the capability to test the depth and pose

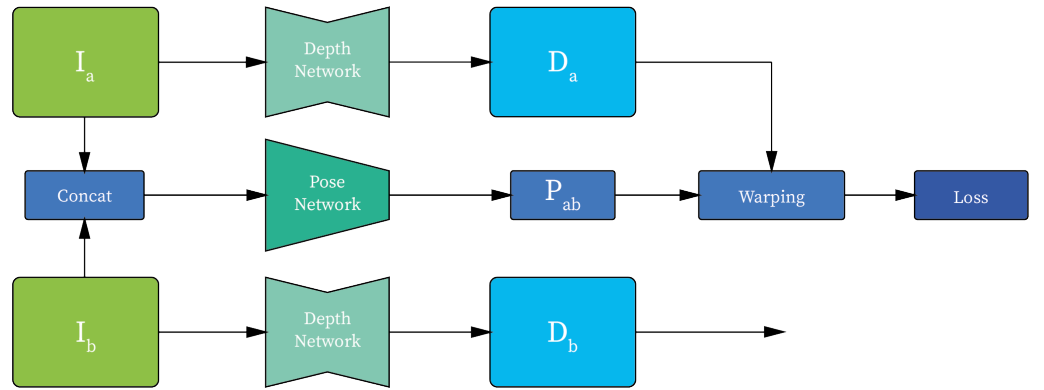


Figure 2. Illustration of the architecture used for the experiments, where I_a , I_b are the input RGB images, D_a , D_b the corresponding estimated depth maps, and P_{ab} is the relative camera position between I_a and I_b .

153 networks independently. We run our first tests on the depth map prediction
 154 on various RGB-D datasets of indoor environments (see Table 2) achieving
 155 results comparable to the ground truth (GT) except for a scale factor that can be
 156 calculated by normalizing the depth map with its median pixel value. The tests
 157 are conducted using the pre-trained model publicly available on the authors'
 158 GitHub repository [26]. We feed the unseen datasets in input to the model and
 159 retrieve the predicted disparity maps. For the evaluation, we adopt the Absolute
 160 Relative difference used in literature which is computed as follow:

$$\frac{1}{|V|} \sum_{p \in V} \frac{|d(p) - d^*(p)|}{d^*(p)} \quad (1)$$

161 where V denotes the set of valid depth pixels, $d(p)$ and $d^*(p)$ are the depth
 162 pixel value of the predicted depth map D and the depth ground truth D^* , respec-
 163 tively. As mentioned before, the predictions are at different scale with respect
 164 to the ground truth. Scaling is then applied via the scaling factor s computed as
 165 follows, where $med\{\}$ refers to the median value:

$$s = \frac{med_{p \in V}\{D^*\}}{med\{D\}} \quad (2)$$

166 Note that, unlike the prediction, the ground truth exhibits some pixels equal
 167 to zero or one due to reflective surfaces or distances out of the sensor range. Such
 168 non-valid pixels are discarded in the computation above.

169 4.1. Dataset

170 The need of virtually reconstructing environments for autonomous naviga-
 171 tion and/or extended reality applications has increased the availability of indoor
 172 RGB-D data to train more and more data-hungry networks; however, the amount
 173 of data is still limited to few common environments. In this section we present a
 174 brief overview of the datasets used in our experiments. We tested the network
 175 performance on four different datasets containing sequences from several indoor
 176 environments. In particular, for the testing purposes we selected the sequences
 177 *freiburg_360* and *freiburg_pioneer* from RGB-D TUM Dataset [24], all the sequences
 178 from RGB-D 7 Scene [17], the RGB-D Scene dataset from Washington RGB-D

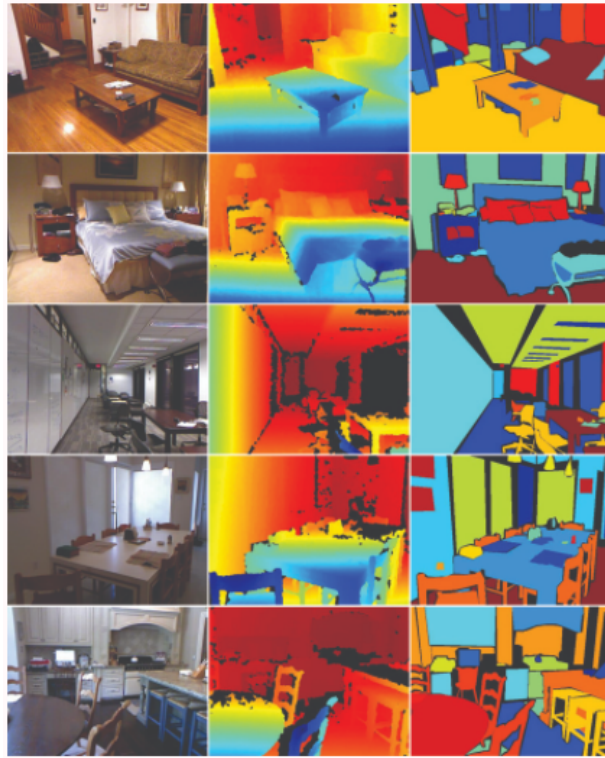


Figure 3. NYU dataset [25]

179 Object Dataset [36] and the SUN RGB-D Dataset [37]. Details about the number
 180 of samples and resolution are reported in Table 2.

- 181 • **RGB-D TUM Dataset:** the sequence *freiburg1_360* contains a 360 degree
 182 acquisition in a typical office environment; the *freiburg_pioneer* sequence
 183 shows a quite open indoor environment captured by a robot with depth
 184 sensor attached on top of it (Figure 4). The dataset is provided with depth
 185 ground truth acquired by the Kinect sensor, and camera pose ground truth
 186 as rotation and translation are acquired with an external motion capture
 187 system, it is typically used for SLAM systems. For additional details we
 188 refer to the dataset website[15] and to the original paper[24]. Among the
 189 available sequences we decided to choose two of them (*freiburg1_360* and
 190 *freiburg_pioneer*) since they represent distinct environments with interesting
 191 characteristics useful to test the generalization of the network. In particular,
 192 in *freiburg_360* there are many complex geometry defined by the office
 193 furniture; *freiburg_pioneer* is instead characterized by wide spaces, usually
 implying more homogeneous depth maps but larger depth range.



Figure 4. RGB-D TUM Dataset, frame taken from the two sequences.

- 194 • **RGB-D Microsoft Dataset:** this dataset [17] consists in sequences of tracked
 195 RGB-D frames of various indoor environments, and it is provided with the
 196 corresponding depth ground truth (Figure 5). This dataset is the one used
 197

198 by the authors in [12] to test the generalization capability of the network.
 199 Accordingly, we decided to re-run the tests as well, to ensure the replicability
 of the paper results.

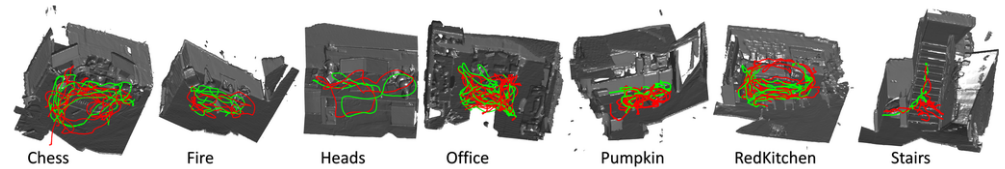


Figure 5. 7 Scene dataset [17]

200
 201 • **Washington RGB-D Object Dataset:** the dataset [36] was created with the
 202 purpose of providing structure data of real objects. Aside the isolated objects,
 203 the dataset provides 22 annotated sequences of various indoor environment
 204 with depth ground truth. Also in this case, RGB-D data are collected using
 Microsoft Kinect using aligned 640x480 RGB and depth images (Figure 6).

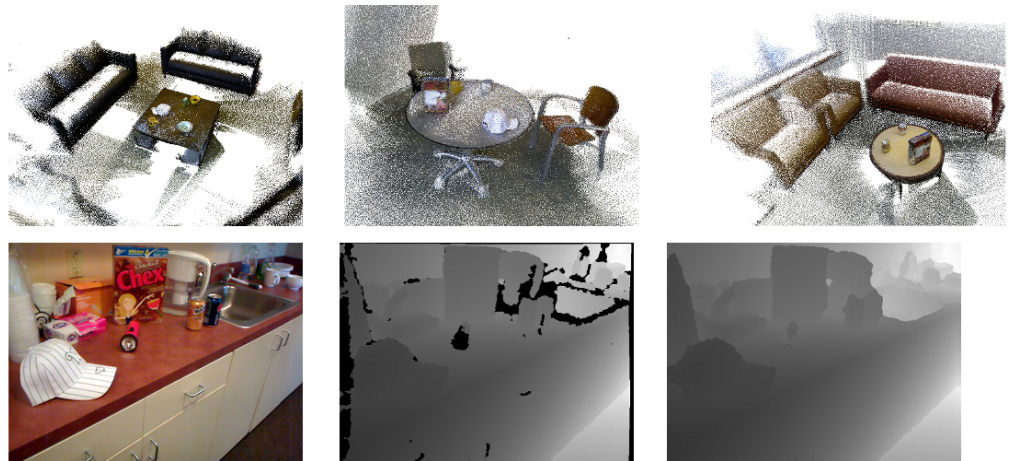


Figure 6. Washington RGB-D Object Dataset [36]

205
 206 • **SUN RGB-D Dataset:** the dataset [37] is a collection of several common
 207 indoor environments from different datasets; it contains RGB-D images from
 208 NYUv2 [27], Berkeley B3DO [38] and SUN3D [39]. The dataset has in total
 209 10335 RGB-D images. In order to make the experiments comparable, we
 have selected only the samples acquired using Kinect (Figure 7).

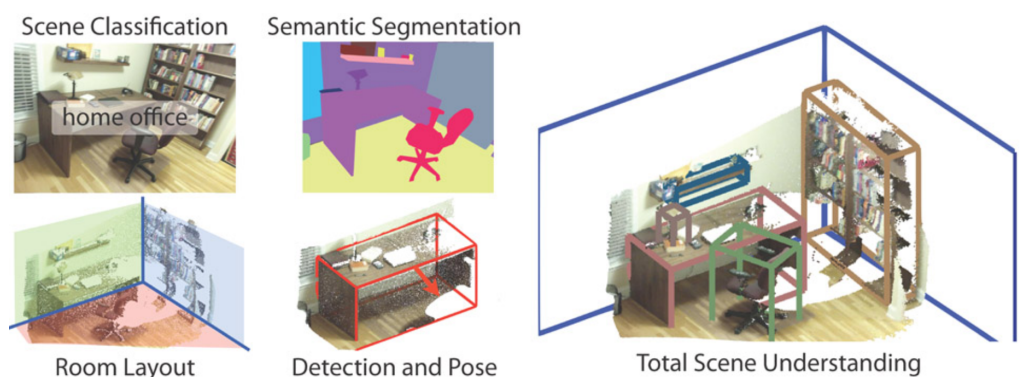


Figure 7. SUN RGB-D Dataset [37]

210

211 As reported above, in all selected datasets, the RGB-D data is acquired with
 212 Microsoft Kinect version 1. The device is equipped with an RGB camera and
 213 a structured light sensor working on the near infrared light spectrum, where a
 214 known infrared pattern is projected onto the scene and the depth is computed
 215 after distortion correction. For additional information about the sensor and the
 216 related performances, please refer to the study by Wasenmüller *et al.* [40]. In
 217 terms of accuracy, the sensor exhibits an exponentially increasing offset going
 218 from 10mm at 0.5m, of up to 40mm at distance of 1.8m. Although the perfor-
 219 mances of the sensor are not as accurate as other more recent devices made
 220 available on the market, most benchmark datasets in the literature still have the
 221 Kinect depth map as ground truth.

Name	#Images	Img. Size	Ref.
<i>freiburg_360</i> (TUM RGB-D)	756	640x480	[24]
<i>freiburg_pioneer</i> (TUM RGB-D)	1225	640x480	[24]
7Scene	29000	640x480	[17]
Washington	11440	640x480	[36]
SUN	10335	640x480	[37]

Table 2: Details of the three dataset used in the testing phase.

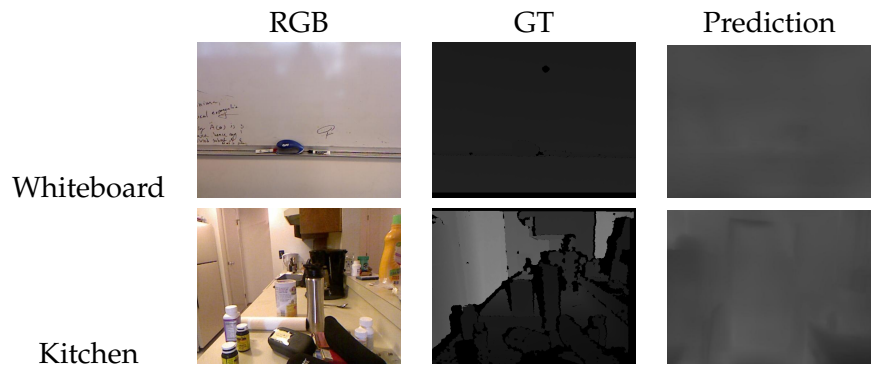


Figure 8. Example of depth map prediction with different depth standard deviation.

222 5. Results

223 In this section we present the results we obtained in our simulations. Since
 224 the author of [12] already compared the network performances with previous
 225 state-of-the-art unsupervised methods, and in particular with [11] and [41] show-
 226 ing an improvement in terms of Absolute Relative error after training data
 227 rectification, we focus on enriching the benchmark by testing the network on
 228 different unseen data. We evaluate the datasets described in the previous section
 229 by feeding frame sequences to the network and computing the Absolute Relative
 230 difference (AbsRel) for each prediction-ground truth pair every 5 frames. The
 231 results are reported in Table 3. We notice that the network generalization perfor-
 232 mance highly depends on the images depth range that has to be estimated. As
 233 an example, environments containing various structural features are more likely
 234 to result in a higher error, and frames depicting an homogeneous scenario with
 235 lower depth variation result in a lower error.

Scenes	AbsRel	StdDev (σ^2)
<i>freiburg_360</i> (TUM RGB-D) [24]	0.16	5056.86
<i>freiburg_pioneer</i> (TUM RGB-D)[24]	0.28	11370.31
<i>Chess</i> (7Scene) [17]	0.19	5800.00
<i>Fire</i> (7Scene) [17]	0.15	4418.00
<i>Office</i> (7Scene) [17]	0.16	4438.00
<i>Pumpkin</i> (7Scene) [17]	0.13	3435.00
<i>RedKitchen</i> (7Scene) [17]	0.20	5700.00
<i>Stairs</i> (7Scene) [17]	0.17	5341.00
Washington [36]	0.30	9656.00
<i>B3DO</i> (SUN RGB-D) [37]	0.18	6886.21

Table 3: Single-view depth estimation results on selected Datasets

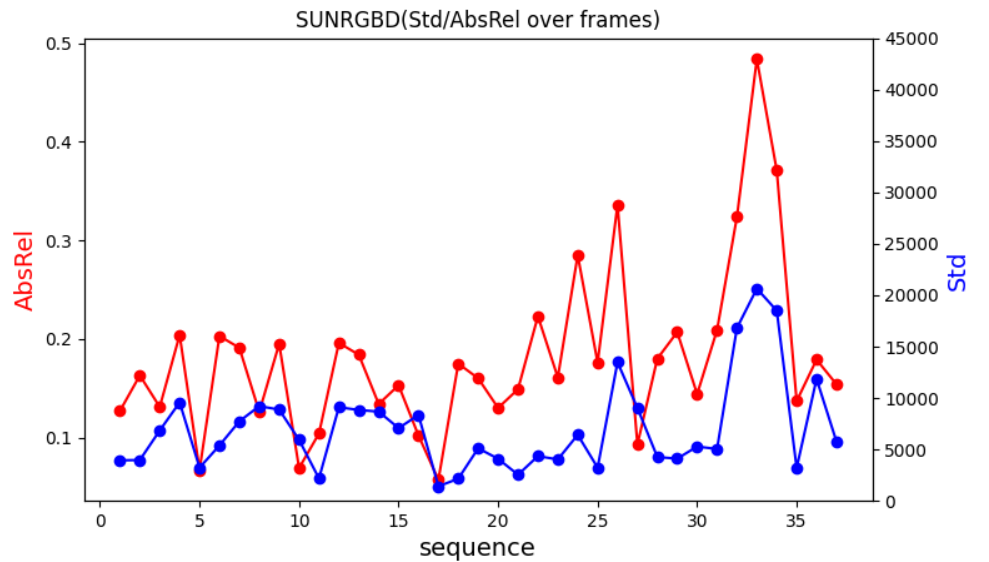


Figure 9. The plot shows the behaviour of Absolute Relative error (AbsRel) and Depth Standard Deviation (Std) over the *B3DO* sequence from SUN RGB-D Dataset.

236 In addition to the Absolute Relative error, we then analysed the Standard
237 Deviation (σ^2) of depth ground truth images, which gives an insight of how
238 challenging an environment is from the learning perspective. The depth standard
239 deviation shows great potential in understanding the overall structure of the
240 environment, thus it can be employed in further improvements of the network
241 depth prediction. As for the AbsRel, the tests were performed computing σ^2
242 along with the error for each frame pairs every 5 frames. Figure 8 shows an
243 example of borderline situations taken from SUN RGB-D [37], where in the
244 case of the whiteboard, the measured AbsRel is particularly low, equal to 0.05
245 and $\sigma^2 = 1416.48$, on the other hand, in the kitchen image the depth range is
246 larger with $\sigma^2 = 20639.78$, and the resulting absolute error is equal to 0.48. By
247 comparing the two examples we can see that frames with a smaller σ^2 consist
248 of relatively simple tasks that the network can easily manage; at the same time
249 they often turn to be *false positives*. This situation is frequent because of the
250 required normalization procedure, which is applied to the predicted depth in
251 order to compare it with the GT. Indeed, for homogeneous surfaces that appear

252 to be orthogonal to the optical axis, the predicted depth map results in an almost
253 flat gray level image, leading, after the normalization, to an apparently optimal
254 prediction, no matter if the scale is consistent or not along the entire sequence.
255 On the other hand, the higher the variation in the depth range, the harder is
256 for the network to predict consistent disparity maps. This behaviour is shown
257 in the plot reported in Figure 9, where the test are conducted on the *B3DO*
258 sequence from SUN RGB-D. Unlike the other sequences, *B3DO* is composed
259 by random frames from different environment, thus it is a good challenge for
260 the generalization capability of the network.. As next step we performed the
261 same test on the remaining (Table 2) to find the contexts in which the network
262 works well and in which ones it is harder for the network to predict the disparity.
263 Figure 10 presents the Absolute Relative error for each considered sequences in
264 relation to the depth Standard Deviation both computed as the mean over the
265 entire sequence. It is arguable from the plot that the Absolute Relative error is
266 directly proportional to the amount of depth information (given by the standard
267 deviation) that the network has to estimate. More precisely, it is noticeable that
268 for datasets such as **7Scene**, **SUN RGB-D** and the sequence *freiburg_360*, where
269 the space is limited and so the overall depth standard deviation, the network
270 tends to remain consistent and more accurate in the prediction, resulting in a
271 lower absolute error. On the other hand, the prediction accuracy decreases when
272 it comes to process wider and more complex environments as the ones belonging
273 to the **Washington** dataset and the sequence *freiburg_pioneer*, and this is due to
274 the higher variation in the environmental depth as it can be seen in Figure 10.

275 6. Conclusions

276 The goal of our paper was to test the generalization performance of the
277 architecture proposed in [11], providing additional benchmark evaluations. The
278 evaluation has been conducted using the Absolute Relative error as a standard
279 metric. In addition we aim at providing the reader with some hints to interpret
280 the reasons behind some of the results we achieved, so as to draw more detailed
281 conclusions. We noticed that the network ability to estimate the structure of an
282 indoor environment is related to the amount of information that has to be learnt,
283 as it can be evinced from the plots reported above. In particular, the data from
284 Washington Dataset shows the worst results and this is mostly due to the larger
285 standard deviation on the depth range. We understand that this parameter can
286 be considered as a valuable parameter to describe the network generalization
287 capability on various environment. According to our experience, we believe the
288 employment of the depth standard deviation as a weighting parameter in the
289 learning stage is a useful parameter to better stimulate the network in predicting
290 consistent disparity maps from large and more complex indoor environments.

291 7. Future Works

292 We tried to extend the evaluation of DispNet in a diversified set of scenarios,
293 with the purpose of testing the depth extraction accuracy in monocular video,
294 using (SoA) CNN. It is needless to say how such an approach can be revolutionary
295 when deployed in real and unconstrained scenarios, and can be proved to be
296 valuable for the companies engaged in the collection of digital twin, as well as

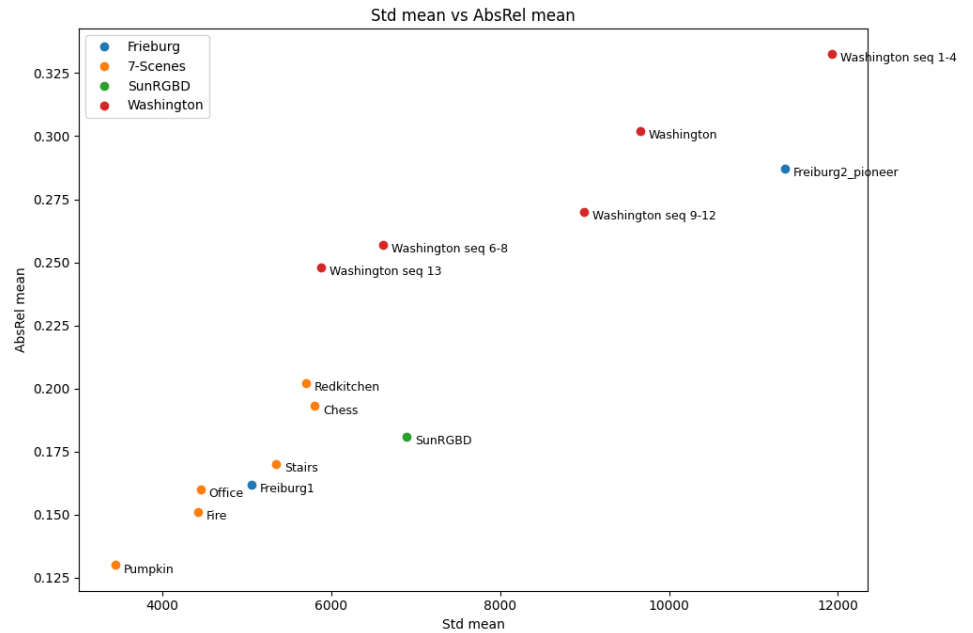


Figure 10. Mean Standard Deviation σ^2 vs. Mean Absolute Relative error of all datasets.

297 for the ones involved in Mixed and Augmented reality developments. Our aim
 298 and recommendation for future studies include:

- 299 • the adoption of other SoA architectures for richer comparisons;
 300 • the adoption of a novel metric that considers the depth standard deviation
 301 for performance evaluation and in the training stage;
 302 • the extension of the study to additional datasets, where the ground truth is
 303 collected with more up-to-date and accurate depth sensors.

304 Abbreviations

305 The following abbreviations are used in this manuscript:

306

SoA	State-of-the-art
SfM	Structure from Motion
SIFT	Scale Invariant Feature Transform
BA	Bundle Adjustemnt
CNN	Convolutional Neural Network
307 DispNet	Disparity Network
RGB	Red, Green, Blue
RGB-D	Red, Green, Blue and Depth
GT	Ground Truth
AbsRel	Absolute Relative error
StdDev	Standard Deviation

References

1. Fazakas, Tamas, and Róbert T. Fekete. "3D reconstruction system for autonomous robot navigation." 2010 11th International Symposium on Computational Intelligence and Informatics (CINTI). IEEE, 2010.
2. Gupta, Sharad Kumar, and Dericks P. Shukla. "Application of drone for landslide mapping, dimension estimation and its 3D reconstruction." *Journal of the Indian Society of Remote Sensing* 46.6 (2018): 903-914.

3. Alexiadis, D. S., D. Zarpalas, and P. Daras. "Real-time, realistic, full 3-D reconstruction of moving humans from multiple Kinect streams." *IEEE Trans. in Multimedia* (2013).
4. Khilar, Rashmita, S. Chitrakala, and SurenderNath SelvamParvathy. "3D image reconstruction: Techniques, applications and challenges." 2013 International Conference on Optical Imaging Sensor and Security (ICOSS). IEEE, 2013.
5. Hosseinian, S., and H. Arefi. "3D RECONSTRUCTION FROM MULTI-VIEW MEDICAL X-RAY IMAGES–REVIEW AND EVALUATION OF EXISTING METHODS." *International Archives of the Photogrammetry, Remote Sensing & Spatial Information Sciences* 40 (2015).
6. Bresnan, Joan, and Sam A. Mchombo. "The lexical integrity principle: Evidence from Bantu." *Natural Language & Linguistic Theory* 13.2 (1995): 181-254.
7. Mayer, Nikolaus, et al. "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.
8. Flynn, John, et al. "Deepstereo: learning to predict new views from real world imagery." U.S. Patent No. 9,916,679. 13 Mar. 2018.
9. Xie, Junyuan, Ross Girshick, and Ali Farhadi. "Deep3d: Fully automatic 2d-to-3d video conversion with deep convolutional neural networks." *European Conference on Computer Vision*. Springer, Cham, 2016.
10. Zhou, Tinghui, et al. "Unsupervised learning of depth and ego-motion from video." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.
11. Bian, Jia-Wang, et al. "Unsupervised scale-consistent depth and ego-motion learning from monocular video." *arXiv preprint arXiv:1908.10553* (2019).
12. Bian, Jia-Wang, et al. "Unsupervised depth learning in challenging indoor video: Weak rectification to rescue." *arXiv preprint arXiv:2006.02708* (2020).
13. Deng, Jia, et al. "Imagenet: A large-scale hierarchical image database." 2009 IEEE conference on computer vision and pattern recognition. Ieee, 2009.
14. He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.
15. Computer Vision Group TUM Department of Informatics Technical University of Munich, RGB-D SLAM Dataset <https://vision.in.tum.de/data/datasets/rgbd-dataset/download>(accessed on May 2021)
16. RGB-D Dataset 7-Scene, <https://www.microsoft.com/en-us/research/project/rgb-d-dataset-7-scenes/>(accessed on April 2021)
17. Glocker, Ben, et al. "Real-time RGB-D camera relocalization." 2013 IEEE International Symposium on Mixed and Augmented Reality (ISMAR). IEEE, 2013.
18. Agrawal, Pulkit, Joao Carreira, and Jitendra Malik. "Learning to see by moving." *Proceedings of the IEEE international conference on computer vision*. 2015.
19. Jayaraman, Dinesh, and Kristen Grauman. "Learning image representations equivariant to ego-motion." *Proc. ICCV*. 2015.
20. Goroshin, Ross, et al. "Unsupervised learning of spatiotemporally coherent metrics." *Proceedings of the IEEE international conference on computer vision*. 2015.
21. Misra, Ishan, C. Lawrence Zitnick, and Martial Hebert. "Shuffle and learn: unsupervised learning using temporal order verification." *European Conference on Computer Vision*. Springer, Cham, 2016.
22. Pathak, Deepak, et al. "Learning features by watching objects move." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017.
23. Wang, Xiaolong, and Abhinav Gupta. "Unsupervised learning of visual representations using videos." *Proceedings of the IEEE international conference on computer vision*. 2015.
24. Sturm, Jürgen, et al. "A benchmark for the evaluation of RGB-D SLAM systems." 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems. IEEE, 2012.
25. NYU depth dataset version 2. https://cs.nyu.edu/~silberman/datasets/nyu_depth_v2.html (accessed on May 2021)
26. <https://github.com/JiawangBian/Unsupervised-Indoor-Depth> (accessed on May 2021)
27. Silberman, Nathan, et al. "Indoor segmentation and support inference from rgb-d images." *European conference on computer vision*. Springer, Berlin, Heidelberg, 2012.
28. Szeliski, Richard. *Computer vision: algorithms and applications*. Springer Science & Business Media, 2010.
29. Wu, Changchang. "Towards linear-time incremental structure from motion." 2013 International Conference on 3D Vision-3DV 2013. IEEE, 2013.
30. Wu, Changchang. "VisualSfM: A visual structure from motion system." (2011).
31. Geiger, Andreas, et al. "Vision meets robotics: The kitti dataset." *The International Journal of Robotics Research* 32.11 (2013): 1231-1237.
32. Lowe, David G. "Object recognition from local scale-invariant features." *Proceedings of the seventh IEEE international conference on computer vision*. Vol. 2. Ieee, 1999.
33. Garg, Ravi, et al. "Unsupervised cnn for single view depth estimation: Geometry to the rescue." *European conference on computer vision*. Springer, Cham, 2016.
34. Vijayanarasimhan, Sudheendra, et al. "Sfm-net: Learning of structure and motion from video." *arXiv preprint arXiv:1704.07804* (2017).
35. Cordts, Marius, et al. "The cityscapes dataset." *CVPR Workshop on the Future of Datasets in Vision*. Vol. 2. 2015.

36. Lai, Kevin, et al. "A large-scale hierarchical multi-view rgb-d object dataset." 2011 IEEE international conference on robotics and automation. IEEE, 2011.
37. Song, Shuran, Samuel P. Lichtenberg, and Jianxiong Xiao. "Sun rgb-d: A rgb-d scene understanding benchmark suite." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.
38. Janoch, Allison, et al. "A category-level 3d object dataset: Putting the kinect to work." Consumer depth cameras for computer vision. Springer, London, 2013. 141-165.
39. Xiao, Jianxiong, Andrew Owens, and Antonio Torralba. "Sun3d: A database of big spaces reconstructed using sfm and object labels." Proceedings of the IEEE international conference on computer vision. 2013.
40. Wasenmüller, Oliver, and Didier Stricker. "Comparison of kinect v1 and v2 depth images in terms of accuracy and precision." Asian Conference on Computer Vision. Springer, Cham, 2016.
41. Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J. Brostow. "Digging into self-supervised monocular depth estimation." In International Conference on Computer Vision (ICCV), pages 3828–3838, 2019.