# RealMonoDepth: Self-Supervised Monocular Depth Estimation for General Scenes

Mertalp Ocal and Armin Mustafa

Center for Vision, Speech and Signal Processing
University of Surrey, UK
{m.ocal,a.mustafa}@surrey.ac.uk

**Abstract.** We present a generalised self-supervised learning approach for monocular estimation of the real depth across scenes with diverse depth ranges from 1–100s of meters. Existing supervised methods for monocular depth estimation require accurate depth measurements for training. This limitation has led to the introduction of self-supervised methods that are trained on stereo image pairs with a fixed camera baseline to estimate disparity which is transformed to depth given known calibration. Self-supervised approaches have demonstrated impressive results but do not generalise to scenes with different depth ranges or camera baselines. In this paper, we introduce RealMonoDepth a self-supervised monocular depth estimation approach which learns to estimate the real scene depth for a diverse range of indoor and outdoor scenes. A novel loss function with respect to the true scene depth based on relative depth scaling and warping is proposed. This allows self-supervised training of a single network with multiple data sets for scenes with diverse depth ranges from both stereo pair and in the wild moving camera data sets. A comprehensive performance evaluation across five benchmark data sets demonstrates that RealMonoDepth provides a single trained network which generalises depth estimation across indoor and outdoor scenes, consistently outperforming previous self-supervised approaches.[1,2]

**Keywords:** Monocular Depth Estimation, Self-Supervised

## 1 Introduction

Humans comprehend 3D from a single viewpoint by leveraging the knowledge of context together with the shape and appearance priors. Similar to human visual perception, robust computer vision systems require the ability to perceive the environment in 3D. This fact has motivated research in monocular depth estimation. The problem to recover depth from a single image is ill-posed due to the projection ambiguity. Supervised methods have been proposed to estimate depth from a monocular image, demonstrating promising results, by training on a large amount of dense ground-truth depth data, however, this

---

[1] Code will be released upon publication.
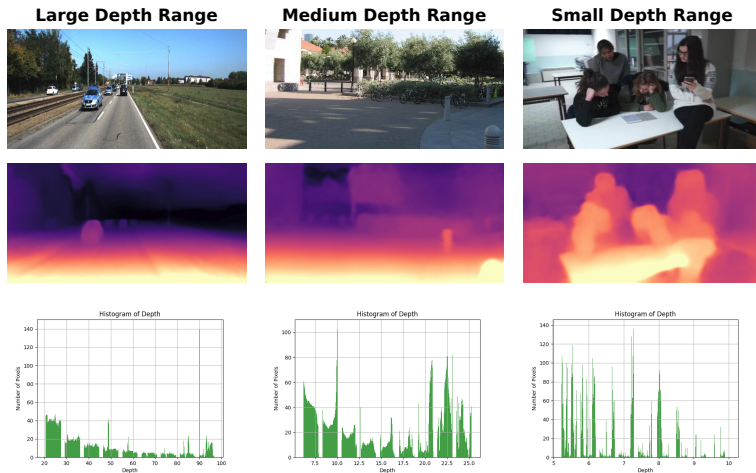
[2] https://youtu.be/6ot3hy3rGaA

**Fig. 1. Top Row:** Scenes with different depth range: left, outdoor, range $\approx 500m$; middle, outdoor, range $\approx 25m$; and right, meeting room, range ($\approx 10m$). **Middle Row:** Depth from proposed method. **Bottom Row:** Depth histogram of ground-truth.

is expensive and impractical to acquire for real-world scenes [12,14,15]. An alternative approach is to generate synthetic data by rendering from computer generated models [41,5,49,58] or 3D scans [8,10,1], but it is challenging to create data that represents the variety and detail of real-world appearance. Also transfer from training on synthetic data to real scenes remains a challenging open problem. Instead of regressing depth from raw pixels, self-supervised learning methods reformulate depth estimation as an image reconstruction problem by re-synthesising a target view from a single source view without ground-truth depth [16,18]. Commonly these methods use stereo image pairs for training with a fixed camera baseline and require scenes with a fixed depth range. These methods learn to estimate inter-image disparity and then use camera calibration to estimate depth.

Learning monocular depth across diverse scenes is a challenging problem, due to large changes in depth range. Typically indoor scenes have a depth range of $< 10m$, whereas outdoor scenes are commonly 100s of meters. Fig. 1 shows some typical scenes with different depth ranges. Monocular depth estimation should be able to estimate depth across scenes with a wide variation in the depth range. It is infeasible to train a deep network that can regress depth values from raw pixels when the output space is incompatible. Existing self-supervised methods can be trained only on data sets with similar depth ranges [19,63,37,62,40], limiting the number of images that can be used for training. As a result, they demonstrate poor generalisation performance and can only perform specific tasks, such as depth estimation in outdoor driving scenes with a fixed stereo baseline. Recently, some supervised approaches [32,33,7] addressed this issue by normalizing the ground-truth to have the same scale which allows learning relative depth values on moving camera data sets with a predefined
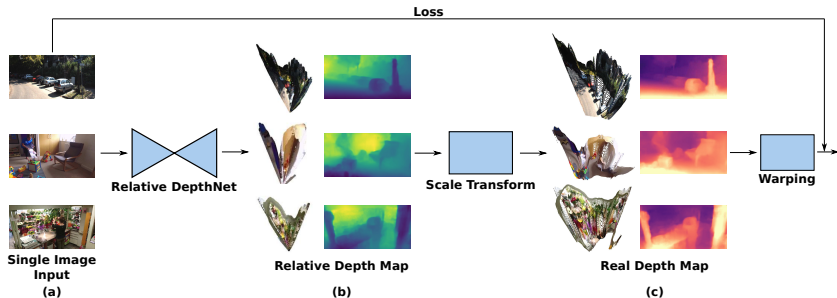
**Fig. 2.** Overview of RealMonoDepth approach for sample input images. For a given single view input image **(a)** the network predicts the **(b)** relative depth map, which is scaled and warped to obtain **(c)** an estimate of the real depth to enable self-supervised training across scenes with diverse depth ranges.

depth range. However, these methods still depend on ground-truth depth values in order to estimate depth from an image.

In this paper, we propose a self-supervised method, RealMonoDepth, that allows a single network to estimate depth for indoor and outdoor scenes demonstrating improved accuracy over previous self-supervised monocular depth estimation approaches. The proposed network learns real depth from stereo-pair and moving camera data sets for scenes with a diverse depth range and without a fixed camera baseline. An overview of our approach for scenes with varying depth ranges (small, medium and large) is depicted in Fig. 1. To enable self-supervised learning of depth estimation across multiple scene scales, we introduce a novel loss function to learn relative depth together with the real depth. The network learns the real depth by transferring relative depth through scaling and warping which is used to compute reconstruction loss for self-supervised training, as shown in Fig. 2. The contributions are of this work are:

- A self-supervised monocular depth estimation method that is able to generalise learning across scenes with different depth range.
- A novel loss function over real depth for self-supervised learning of depth from a single image.
- Evaluation on five benchmark data sets demonstrates generalisation across indoor and outdoor scenes with improved performance over previous work.

## 2   Related Work

Although it is not possible to perform metric reconstruction from a single image due to the projective scale ambiguity, recently proposed learning-based methods have demonstrated that a reliable estimate of the scene geometry can be generated by having prior knowledge of the scale of objects in the scene. This section reviews the supervised, unsupervised and self-supervised approaches that take a single RGB image as input and estimate per-pixel depth as output.

## 2.1    Supervised Depth Estimation

Supervised single image dense depth prediction methods exploit depth values obtained from active sensors such as Kinect and LIDAR as ground-truth for training. Eigen et al. [14] and Fu et al. [15] exploited regression loss to estimate depth for indoor and outdoor scenes respectively. Mayer et al. [38] demonstrated that training a fully convolutional network [36] is an effective approach to learn disparity from stereo images. They also generated a large synthetic data set to train their network called DispNet, demonstrating improved performance on scenes with diverse depth ranges. However the network trained on synthetic data gives limited performance on real in the wild scenes. Huang et al. [23] proposed a deep CNN that leverages multi-view stereo images with known camera poses and calibration to generate patch plane sweep volumes with respect to a reference view. Matching is performed between these patches and produces the inverse disparity map for the reference view. Instead of generating plane sweep volumes, Yao et al. [60] applied differential homography warping on 2D feature maps learned from multiple input views to produce a 3D cost volume. However, these methods suffer from the limited availability of multi-view data of real scenes with ground-truth depth, thereby relying heavily on synthetic data to train their network, which leads poor generalisation capability on complex real scenes. Hence these supervised methods require ground-truth depth or synthetic data for monocular depth estimation on real scenes.

## 2.2    Unsupervised and Self-Supervised Depth Estimation

Recently, unsupervised and self-supervised methods for depth estimation have gained attention eliminating the requirement of ground-truth depth data for real scenes. Unsupervised methods simultaneously estimate the depth and pose from a single RGB image and self-supervised methods exploit camera pose information estimated as a pre-process to obtain depth from a single monocular image.

Zhou et al. [63] pioneered the work in unsupervised depth estimation by proposing separate deep CNN networks for pose estimation between unlabelled video sequences and for single view depth estimation. Instead of training an additional pose estimator network, Wang et al. [54] implemented a differentiable version of Direct Visual Odometry Space (DVO) [51] which is popularly used in current SLAM [39,11] algorithms. Specifically, DVO solves for the pose by minimizing the warping loss of the reference frame from the source frame given the reference frame depth. Furthermore, they introduced a depth normalisation layer in order to address the scale ambiguity problem, which significantly outperformed [63]. Inspired by the relation between depth, pose and optical flow tasks in 3D scene geometry, Yin et al. [61] jointly trained an unsupervised end-to-end deep network to predict pose and depth for non-rigid objects. An optical flow consistency check is imposed between backward and forward flow estimations for reliable estimation. However all these unsupervised methods give a limited performance on general scenes because of the ambiguity in the projection scale introduced by both unknown depth and pose.

Self-supervised methods exploit known camera pose information to resolve the depth ambiguity given estimates of disparity. Garg et al. [16] introduced a self-supervised learning approach to train a deep network from stereo pairs by exploiting the epipolar relation between the cameras given a known calibration in order to generate inverse warp of the left view to reconstruct the right view. Although their results are impressive, they use non-differentiable Taylor series expansion to perform warping of disparity. Godard et al. [18,19] imposed left-right consistency as a constraint for disparity regularisation and established a differentiable optimisation by leveraging spatial transformer networks [24] for bilinear sampling. Given a stereo pair as input for training, they estimate two disparity maps: left view disparity with respect to right view and right view disparity with respect to left view. Then, they reconstruct both views and also the disparity maps to achieve left-right consistency. Moreover, they also imposed edge-aware smoothness as another regulation along with left-right consistency and leveraged SSIM [56] loss in addition to L1 reconstruction loss. Poggi et al. [42] proposed an improvement to [18] by leveraging three rectified views instead of stereo in order to establish additional disparity consistency. Inspired by recent successes of deep learning in single image super-resolution, Pillai et al. [40] employ sub-pixel convolutional layers instead of resizing convolution layers. Contrary to previous methods that are limited to low-resolution operation, their method exploits high fidelity for better self-supervision. Existing self-supervised approaches estimate disparity and assume a fixed camera baseline during training. This limits the approaches to training for scenes with a similar depth range and does not allow generalisation across diverse indoor and outdoor scenes, or diverse data sets for training. In this paper, we introduce an approach to self-supervised learning using a loss function based on estimates of the true scene depth. This allows generalisation across both stereo pair and moving camera data sets for scenes with different depth ranges.

Both self-supervised and unsupervised methods suffer from the following limitations: requirement of fixed camera baseline; limited generalisation performance on scenes with varying depth ranges (indoor/outdoor); and self-supervised methods estimate disparity and only work for training on stereo image pairs with a fixed baseline which limits the training set. The proposed self-supervised depth estimation method addresses all of these limitations by generalising learning across scenes with different depth ranges and works for both stereo and moving camera data sets improving generalisation and accuracy of depth estimation over previous approaches.

## 3   Method

This paper introduces a self-supervised single image depth estimation approach that is able to generalise learning across scenes with diverse depth ranges. The method is trained on both stereo image pairs and moving camera data sets, giving state-of-the-art performance across five benchmark data sets. The method estimates depth from a single view image, an overview is shown in Fig. 2. The
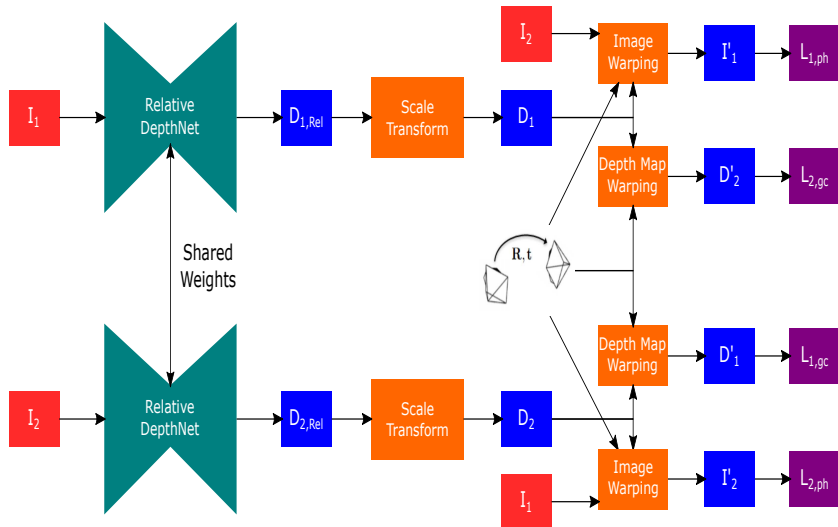
**Fig. 3.** Training framework for proposed self-supervised single image depth estimation

proposed network is trained on two views of a static scene. The Relative Depth-
Net network estimates relative depth maps from two views, inspired from [19]
which estimates disparity maps between stereo pair. As a pre-process, camera
calibration is estimated from sparse correspondences between two views using
an existing visual SFM method COLMAP [48]. The sparse correspondences are
used to estimate the median scene depth for scale transformation, which is ap-
plied on the relative depth maps to obtain real depth maps/true scene depth.
The real depth maps and images are warped between views using the calibration
to estimate the loss. This enables the network to be trained on moving camera
and stereo data sets of real scenes with varying depth ranges and camera view
baselines to generalise performance across a variety of indoor and outdoor scenes
in the wild.

### 3.1   Training Framework

The training framework for the proposed approach is illustrated in Fig. 3 for
learning single view depth estimation from two viewpoint images. Given two
images of a scene from different viewpoints $(I_1, I_2)$, the depth network (Relative
DepthNet) predicts the corresponding per-pixel relative depth maps $(D_{1,Rel}, D_{2,Rel})$ using shared weights. The relative depth maps are transformed to real
depth $(D_1, D_2)$ using the scale transform module. The self-supervised loss is
then computed using the warped real depth estimates $(D_1, D_2)$ and warped
images $(I_2, I_1)$. Using estimated calibration and camera poses, real depth values
allow us to reconstruct the input images $(I'_1, I'_2)$ and depth maps $(D'_1, D'_2)$. This
information is interpolated to compute photometric loss $(L_{ph})$ and geometric
consistency $(L_{gc})$ loss, that supervises the depth network. SSIM and smoothness
losses are introduced to regulate the depth estimation. The loss function for the
proposed method is:

$$L_{final} = \sum_s \lambda_{ph} L_{ph}^s + \lambda_{gc} L_{gc}^s + \lambda_{ssim} L_{ssim}^s + \lambda_s^s L_{smooth}^s \qquad (1)$$

where $s$ indexes over different image scales and $\lambda_{ph}$, $\lambda_{gc}$, $\lambda_{ssim}$ and $\lambda_{smooth}$ are the weighting terms. Extension to training with multiple viewpoint images ($\wr 2$) is straightforward (using loss $L_{i,ph}, L_{i,gc}$ based on the estimated real scene depth $D_i$ for each input image $i$). While existing self-supervised methods are only suited for training on fixed baseline cameras and scenes with similar depth range, we overcome this limitation by stabilizing the output space of the depth network, allowing training on indoor and outdoor scenes with diverse depth ranges.

**Relative DepthNet Architecture:**
Based on the U-Net Architecture [44], in order to effectively capture both local and global information, we use a multi-scale encoder-decoder network with skip connections, similar to [18]. Inspired by [21,19,27], we select ResNet50 [22] for the encoder with initialised weights pre-trained on ImageNet [46] and randomly initialise decoder weights. Unlike popular self-supervised/unsupervised single depth estimation approaches [18,19,37,63,40,31], our network predicts relative depth instead of inverse depth or disparity. We observe that applying sigmoid activation at the network output slows convergence for close and far depth values due to the vanishing gradient problem. Instead, we replace sigmoid with identity activation and handle negative values with exponential mapping, detailed in the Scale Transform section. Apart from multi-scale output layers of the decoder, we apply batch normalisation and ReLU nonlinearities in all layers. A detailed explanation with the full details of the network architecture required for implementation is presented in the supplementary material.

**Scale Transform:**
To learn to estimate depth from images across diverse scenes with varying depth ranges, we normalise depth across the images and data sets using a non-linear scale transform and train the network to estimate relative depth. Given the relative depth map prediction as input, our scale transform module outputs the real depth map. This is formulated as:

$$D_k = \mu_k e^{(D_{k,Rel})} \text{ for } k = 1, 2 \qquad (2)$$

where $\mu_k$ is the median depth value for two images $I_1$ and $I_2$, $D_1$ and $D_2$ are the real scene depth maps. Inspired by supervised methods [14,15], we use exponential mapping in equation 2 in order to reduce the penalisation of the deep network from distant and ambiguous depth values.

During training, camera calibration is required to estimate the median depth value for the scale transform. For data sets with unknown calibration (e.g. Mannequin data set), off-the-shelf SFM method COLMAP [48] is used to solve for camera calibration and sparse correspondences between views. If the camera calibration is known (e.g. KITTI data set), we use the calibration to compute sparse correspondences between views. The sparse correspondences are triangulated in 3D exploiting camera pose to get a sparse reconstruction of the scenes.

The sparse 3D points are projected on each view to obtain sparse depth maps for each viewpoint. Depth values are then sorted and the median depth value is estimated, as shown in Fig. 1. Median depth enables prediction of real depth maps which are used together with the input images to estimate the loss in Equation 1.

**Loss Functions:**
**Photometric consistency:** This loss enables estimation of depth and is inspired by self-supervised learning approaches that reformulate depth estimation as an image reconstruction problem [16,18]. Here, the underlying intuition is that every viewpoint is a 2D projection of the same 3D scene, so one view can be reconstructed from another view, which implies knowledge about depth. With known calibration between views, depth is treated as an intermediate variable to perform novel-view synthesis with a deep network. Here, $I'_1$ is the reconstructed reference view from $I_2$ and $I'_2$ is the reconstructed reference view from $I_1$. Let $K_1$, $K_2$ matrices represent the intrinsic parameters of $I_1$ and $I_2$ respectively, $P_{1 \to 2}$, $P_{2 \to 1}$ denote relative camera pose matrices between views, $D'_2(p'^{(2)})$, $D'_1(p'^{(1)})$ represent projected depth values for each pixel $p$. Then, the homogeneous pixel-wise projection relations $(D'_{proj}())$ between two views can be formulated as:

$$D'_{proj,2}(p'^{(2)})p'^{(2)} = K_2 P_{1 \to 2} D_1(p^{(1)}) K_1^{-1} p^{(1)},$$
$$D'_{proj,1}(p'^{(1)})p'^{(1)} = K_1 P_{2 \to 1} D_2(p^{(2)}) K_2^{-1} p^{(2)} \tag{3}$$

Since the projected pixel values are continuous, we use a spatial transformer network [24] to perform differentiable bilinear sampling in order to approximate $I'_1$ and $I'_2$ by interpolating $I_1(p^{(1)})$ and $I_2(p^{(2)})$ from neighboring corner pixels:

$$I'_1(p^{(1)}) = I_2(p^{(2)}) = \sum_{(i,j)} w^{(2)}_{(i,j)} I_2(p'^{(2)}_{(i,j)}) \ , \ I'_2(p^{(2)}) = I_1(p^{(1)}) = \sum_{(i,j)} w^{(1)}_{(i,j)} I_1(p'^{(1)}_{(i,j)})$$

Here, $(i, j)$ are the indices of the top left, top right, bottom left and bottom right pixels of the projected pixels and $w_{(i,j)}$ is the corresponding weight, which is inversely proportional to spatial distance. Also, due to occlusions between viewpoints, some pixels in the reference view will be projected outside of the image plane boundary in the source view. In order to prevent these unresolved regions from penalising the network, similar to Mahjourian et al. [37], a validity binary masks is computed to exclude these pixels from the training loss [19]. Hence, the L1 image reconstruction loss for two views $k = 1, 2$ becomes:

$$L_{ph} = \frac{1}{N} \sum_{p}^{N} \sum_{k=1,2} L_{k,ph}(p) \text{ where, } L_{k,ph}(p) = |I_k(p) - I'_k(p)| M_k(p) \tag{4}$$

Also, SSIM [56] loss is applied to 3x3 image patches in order to regulate the noisy artifacts caused by the L1 loss, defined as follows:

$$L_{ssim} = \frac{1}{N} \sum_{p}^{N} [(1 - SSIM(I_1, I'_1)(p)) + (1 - SSIM(I_2, I'_2)(p))] \tag{5}$$

**Smoothness:** This loss term ensures the depth maps are smooth, reducing the noise in the depth maps. Sobel gradients are used to calculate the loss instead of horizontal and vertical gradients in Wong et al. [57]. Sobel gradients allow depth to be smooth horizontally, vertically and diagonally. In most cases, regions that have higher reconstruction error are only visible in one view. Applying smoothness regularisation to these unresolvable regions induces a false penalty in the network training. In order to overcome this problem, an adaptive smoothness regulation weight is used for every pixel that varies in space, as in previous work [57]. Based on Equation 4, the adaptive weight for a pixel is computed as follows: $\alpha_k(p) = exp(-\frac{c*L_{k,ph}(p)}{\sigma_k})$, where index $k = 1, 2$ corresponds to the view and $c$ is the scale factor that determines the range of $\alpha$ and $\sigma$ is the global residual represented as: $\sigma_k = \frac{1}{\frac{1}{N}\sum\limits_{p}^{N} L_{k,ph}(p)}$. $\alpha$ depends on global residual ($\sigma$) at each position and tends to be small when the residuals are high. The average of $\alpha$ approaches to 1 as the training converges. Based on adaptive weight, the smoothness regularisation objective for two views $k = 1, 2$ is:

$$L_{smooth} = \frac{1}{N}\sum_{p}^{N}\sum_{k=1,2} L_{k,smooth}(p) \qquad (6)$$

$$L_{k,smooth}(p) = \alpha_k(p)(|\partial_x D_{k,rel}(p)|exp(-\partial_x I_k(p)) + |\partial_y D_{k,rel}(p)|exp(-\partial_y I_k(p))$$

Here, the depth values are enforced to be locally smooth which is represented as the x and y Sobel gradients of the relative depth maps.

**Geometric consistency:** The per-view predicted depth maps may not be consistent with the same 3D geometry. This causes depth discontinuities and outliers on the surfaces and boundaries of the objects. In order to address this issue, we learn the real depth map ($D_1$, $D_2$) for each viewpoint simultaneously with consistency checks. Inspired by Bian et al. [3], we enforce geometric consistency symmetrically by sampling different viewpoints in the same training batch. Based on Equation 3, we interpolate projected depth maps $D'_{proj,1}$, $D'_{proj,1}$ with bilinear sampling in order to approximate $D'_1(p)$, $D'_2(p)$ which lie on the pixel grid. The geometric consistency loss function for $k = 1, 2$ is defined as follows:

$$L_{gc} = \frac{1}{N}\sum_{p}^{N}(L_{1,gc}(p) + L_{2,gc}(p)) \text{ where } L_{k,gc}(p) = \frac{|D_k(p) - D'_k(p)|}{D_k(p) + D'_k(p)} \qquad (7)$$

Here, similar to [3] we use a normalised symmetric loss function to achieve depth consistency between predicted depth maps for each viewpoint.

## 4   Experiments

Qualitative and quantitative results are presented on five benchmark data sets against state-of-the-art supervised, unsupervised and self-supervised methods.

**Table 1.** Data sets for monocular depth estimation

| data set | Indoor | Outdoor | Dynamic | Video | Depth | Diversity | Annotation | # Images |
|---|---|---|---|---|---|---|---|---|
| NYUDv2 [50] | ✓ | | ✓ | ✓ | Metric | Low | RGB-D | 407K |
| Make3D [47] | | ✓ | | | Metric | Low | Laser | 534 |
| KITTI [17] | | ✓ | ✓ | ✓ | Metric | Low | Stereo | 93K |
| DIW [7] | ✓ | ✓ | ✓ | | Ordinal Pair | High | User clicks | 496K |
| Cityscapes [9] | | ✓ | ✓ | ✓ | Metric | Low | Stereo | 25K |
| Megadepth [33] | | ✓ | ✓ | | No scale | Medium | SFM | 130K |
| MC [32] | ✓ | ✓ | ✓ | ✓ | No scale | High | SFM | 115K |
| TUM [52] | ✓ | | ✓ | ✓ | Metric | Low | RGB-D | 80K |

We demonstrate that the proposed self-supervised loss function using real depth dramatically improves generalisation performance when trained on both moving camera (Mannequin Challenge (MC) [32] mostly indoor) and stereo (KITTI [17] outdoor) data sets jointly. These data sets contain both indoor (1–10m) and outdoor (1–1000m) scenes with a wide variation in depth range. We test the same trained model on four benchmark data sets which the network has not seen during training: KITTI Eigen test split [13] (street scenes), Make3D [47] (outdoor buildings), NYUDv2 test split [50] (indoor) and dynamic subset of TUM-RGBD [32] (humans in indoor environments). Key attributes of the data sets used in experiments are listed in Table 1. Additional qualitative results on a wide variety of in the wild scene images for the DIW [7] data set are presented in the supplementary material, together with comparative performance evaluation on a diverse range of challenging in the wild videos in the supplementary video.

## 4.1   Implementation Details

Our model is implemented in Tensorflow [2], trained using the Adam [26] opti-miser for 25 epochs with an input/output resolution of $512 \times 256$ and a batch size of 12. Each batch sample consists of different viewpoint images of the same scene. We set the number of viewpoint images as 2 which leads to $12 \times 2 = 24$ images for each batch. Initial learning rate is set to $10^{-4}$ and it is decayed with a lin-ear scheduler $LR(iteration) = initialLR * (1 - iteration/maxIteration)^{0.9}$. The weights for the loss terms are empirically determined as: $\lambda_{ph} = 0.15, \lambda_{ssim} = 0.85, \lambda^s_{smoothness} = 0.01/s$ where $s$ is the downsampling factor for each scale, $\lambda_{gc} = 0.1$ and scale factor for adaptive regularisation term is chosen as 5, similar to [57]. For data augmentation, we perform horizontal flipping, random scal-ing, cropping with 50% and apply random brightness, contrast, saturation on the rest 50%, with the same range of values as in [19]. Full details of network implementation are given in the supplementary material.

## 4.2   Experimental Setup

We train our model based on data split of Eigen et al. [14] for KITTI and Li et al. [32] for Mannequin Challenge (MC) data sets. We perform both individual and mixed training on these data sets with and without our proposed scale transform approach in order to evaluate the differences in generalisation capability and

compare with previous methods. The underlying motivation for combining these data sets is threefold: 1) They are both large data sets suitable for self-supervised training, 2) Their sizes are comparable which makes them favourable to mix for joint training, 3) They represent different varieties of appearance: outdoor street scenes in KITTI and humans (mostly indoors) in MC.

We select 23,488 stereo pairs of KITTI for training and the remaining 697 images are used as test set for evaluating single view depth estimation. Similar to train split of [32], we select 2463 scenes of MC for training. Since some of the videos on Youtube were deleted by the owners, we were not able to access all of the video URLs provided by [32], so our training set is slightly smaller. In order to ensure well balanced class distribution, we randomly sample 40 viewpoint images from each scene. Images are resampled for scenes that have a lower number of viewpoints.

We quantitatively evaluate the single view depth estimation of our model following the error metrics of Eigen et al. [14]: mean absolute relative error (Abs Rel), mean squared relative error (Sq Rel), root mean squared error (RMS), and root mean squared log10 error (RMS(log)). Following Zhou et al. [63], we scale our relative single-view depth map predictions to match the median of ground-truth.

## 4.3   Comparison with State-of-the-art

This section presents quantitative results of the proposed scale transform method models trained on various data set combinations (KITTI, MC, MC+KITTI). Note: due to estimation of disparity rather than depth and assumption of a fixed camera baseline in previous self-supervised estimation methods [18,20,42,16], it is not possible to train for data sets with different scene scales.
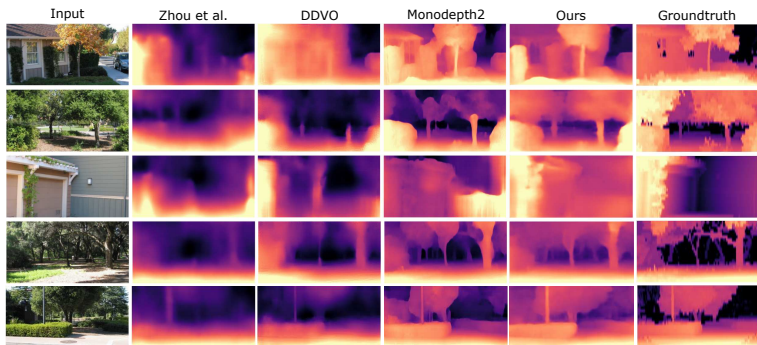
**KITTI** We report the results for the test set of KITTI based on Eigen split [14] in Table 2. Our model trained on MC+KITTI shows the best generalisation performance outperforming all state-of-the-art supervised methods which are not trained on KITTI and which are trained on large scale diverse data sets such as Chen et al. [6] and Li et al. [33]. We also demonstrate that the training loss function based on real depth allows generalisation over data sets for indoor and outdoor scenes with different scales and allows the combination of data sets during training without degrading test performance. Our model trained on MC+KITTI outperforms other state-of-the-art self-supervised/unsupervised methods on KITTI even when they are trained on KITTI.

**Make3D** Next, we evaluate on the Make3D data set following the procedure in [19]. In Table 3, our model trained on MC+KITTI demonstrates the best generalisation performance on an unseen data set compared to other methods which are [35,34] and are not [14,7] trained on Maked3D. Qualitative comparisons of our method trained on MC+KITTI with training loss based on real depth demonstrates improved performance in Fig.4.

**NYUDv2** Here, we show that our proposed method also generalises well to indoor scenes. We evaluate performance against state-of-the-art self-supervised monocular method Monodepth2[19] on the NYUDv2 test split. In Table 3, our

**Table 2.** Quantitative results on KITTI data set for different methods trained on various scenes (lower is better).

| Method | Supervision | Training set | Abs Rel | Sq Rel | RMS | RMS(log) |
|---|---|---|---|---|---|---|
| Eigen [14] | Depth | KITTI | 0.203 | 1.548 | 6.31 | 0.282 |
| Liu [34] | Depth | KITTI | 0.202 | 1.614 | 6.52 | 0.275 |
| Fu [15] | Depth | KITTI | **0.072** | **0.307** | **2.73** | **0.120** |
| Eigen [14] | Depth | NYU | 0.521 | 5.016 | 10.37 | 0.510 |
| Liu [34] | Depth | NYU | 0.540 | 5.059 | 10.10 | 0.526 |
| Laina [28] | Depth | NYU | 0.515 | 5.049 | 10.07 | 0.527 |
| Liu [34] | Depth | Make3D | 0.362 | 3.465 | 8.70 | 0.447 |
| Laina [28] | Depth | Make3D | 0.339 | 3.136 | 8.68 | 0.422 |
| Chen [7] | Depth | DIW | 0.393 | 3.260 | 7.12 | 0.474 |
| Li [33] | Depth | Megadepth | 0.368 | 2.587 | 6.68 | 0.414 |
| Garg [16] | Pose | KITTI | 0.152 | 1.226 | 5.85 | 0.246 |
| Monodepth [18] | Pose | KITTI | 0.148 | 1.334 | 5.93 | 0.247 |
| DDVO [54] | | KITTI | 0.151 | 1.257 | 5.58 | 0.228 |
| GeoNet [61] | | KITTI | 0.149 | 1.060 | 5.57 | 0.226 |
| Struct2depth [4] | | KITTI | 0.141 | 1.026 | 5.29 | 0.215 |
| Zhou [63] | | KITTI | 0.208 | 1.768 | 6.86 | 0.283 |
| Zhou [63] | | CS | 0.267 | 2.686 | 7.58 | 0.334 |
| 3Net [42] | Pose | KITTI | 0.129 | 0.996 | 5.28 | 0.223 |
| Monodepth2 [19](640x192) | Pose | KITTI | 0.109 | 0.873 | **4.960** | 0.209 |
| Ours | Pose | KITTI | 0.109 | 0.928 | 4.99 | 0.199 |
| Ours | Pose | MC | 0.276 | 2.563 | 9.17 | 0.386 |
| **Ours** | Pose | MC+KITTI | **0.108** | **0.855** | 5.15 | **0.204** |



**Fig. 4.** Qualitative results on the Make3D data set.

model outperforms Monodepth2 with a significant margin and achieves competitive accuracy against supervised methods that are trained on a different split of the same data set. We also provide qualitative comparisons of our model trained on MC+KITTI in Fig. 5 demonstrating improved performance.
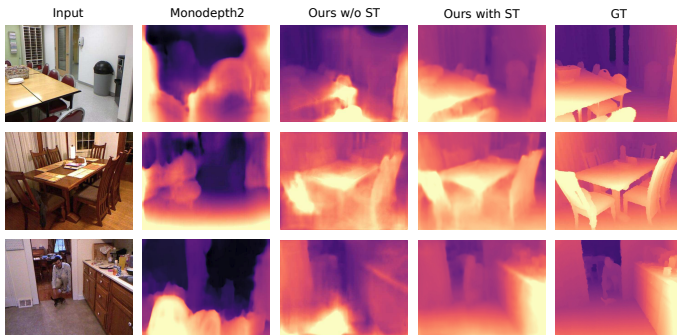
**TUM** Finally, we present quantitative and qualitative results tested on the dynamic subset of the TUM-RGBD data set in Table 4 and Fig. 6 respectively. Our model ranks second-best compared to supervised methods and best for self-supervised methods. In order to make a fair comparison with [32], we only include their result trained on a single image without any additional prior knowledge except the depth ground-truth.

## 4.4   Ablation Study

Here, we evaluate the effect on generalisation performance across scenes with different depth ranges for our self-supervised training using a loss function based

**Table 3.** Quantitative results on Make3D and NYUDv2 data sets for different methods trained on various scenes (lower is better).

| | Make3D | | | | | NYUDv2 | | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | Supervision | Training set | Abs Rel | RMS | Method | Supervision | Training set | Abs Rel | RMS |
| Xu [59] | Depth | Make3D | 0.184 | **4.38** | Xu [59] | Depth | NYU | **0.121** | 0.586 |
| Li [30] | Depth | Make3D | 0.278 | 7.19 | Li [29] | Depth | NYU | 0.139 | **0.505** |
| Laina [28] | Depth | Make3D | **0.176** | 4.45 | Laina [28] | Depth | NYU | 0.129 | 0.583 |
| Liu [34] | Depth | Make3D | 0.314 | 8.60 | Liu [34] | Depth | NYU | 0.230 | 0.824 |
| Liu [35] | Depth | Make3D | 0.335 | 9.49 | Liu [35] | Depth | NYU | 0.335 | 1.06 |
| Laina [28] | Depth | NYU | 0.669 | 7.31 | Eigen [14] | Depth | NYU | 0.215 | 0.907 |
| Liu [34] | Depth | NYU | 0.669 | 7.20 | Eigen [13] | Depth | NYU | 0.158 | 0.641 |
| Eigen [14] | Depth | NYU | 0.505 | 6.89 | Roy [45] | Depth | NYU | 0.187 | 0.744 |
| Chen [7] | Depth | DIW | 0.550 | 7.25 | Wang [55] | Depth | NYU | 0.220 | 0.745 |
| Li [33] | Depth | Megadepth | 0.402 | 6.23 | Jafari [25] | Depth | NYU | 0.157 | 0.673 |
| Monodepth [18] | Pose | KITTI | 0.525 | 9.88 | Monodepth2 [19] | Pose | KITTI | 0.342 | 1.183 |
| Monodepth2 [19] | Pose | KITTI | 0.322 | 7.42 | Ours | Pose | KITTI | 0.300 | 1.005 |
| Zhou [63] | | KITTI | 0.651 | 8.39 | Ours | Pose | MC | 0.201 | 0.718 |
| DDVO [54] | | KITTI | 0.387 | 8.09 | **Ours** | Pose | MC + KITTI | **0.193** | **0.686** |
| Ours | Pose | KITTI | 0.295 | 7.10 | | | | | |
| Ours | Pose | MC | 0.346 | 7.70 | | | | | |
| **Ours** | Pose | MC+KITTI | **0.289** | **6.92** | | | | | |



**Fig. 5.** Qualitative results on NYU data set.

**Table 4.** Quantitative results on TUM Dynamic Objects RGBD data set for different methods trained on various scenes. Lower is better.

| Method | Supervision | Training set | Abs Rel | RMS |
|---|---|---|---|---|
| Xu [43] | Depth | NYU | 0.274 | 1.085 |
| Laina [28] | Depth | NYU | 0.223 | 0.947 |
| Chen [7] | Depth | NYU+DIW | 0.262 | 1.004 |
| DeMoN [53] | Depth | TUM RGBD+MVS | 0.220 | 0.866 |
| Fu [15] | Depth | NYU | **0.194** | 0.925 |
| Li [32](single image) | Depth | MC | 0.204 | **0.840** |
| Monodepth2 [19] | Pose | KITTI | 0.427 | 1.616 |
| Ours | Pose | KITTI | 0.388 | 1.533 |
| Ours | Pose | MC | 0.207 | 0.973 |
| **Ours** | Pose | MC+KITTI | **0.201** | 1.025 |

on real depth estimates. For models trained without our proposed method, we omit the scene median depth term, so the scale transform module is modified as $D = exp(D_{rel})$ where $D_{rel}$ is the estimated relative depth of our Relative DepthNet network and $D$ is the real depth map which is used for warping the other viewpoint image in order to compute the loss for training. We train our network with four different configurations regarding training data sets and usage of proposed scale transform (ST) during training: 1) MC w/o ST, 2) MC+KITTI
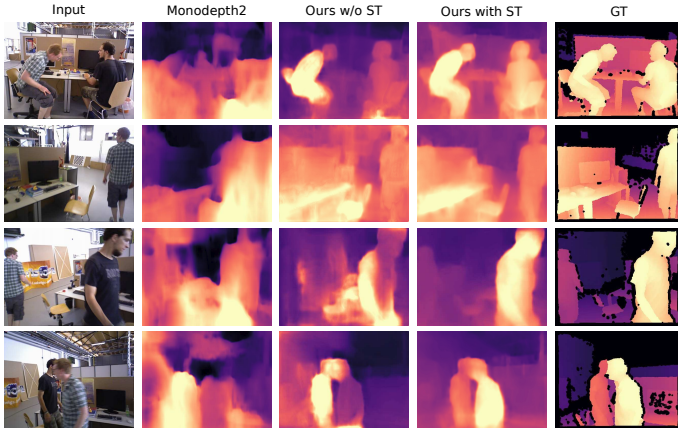
**Fig. 6.** Qualitative results on TUM data set.

w/o ST, 3) MC w/ ST and 4) MC+KITTI w/ ST. Then, we test each of these four models on four benchmark data sets similar to section 4.3. We report numerical ablation results in Table 5 and show qualitative comparisons between models trained on MC+KITTI with and without scale transform in Fig. 5, Fig. 6 and Fig. 7. Both our models trained on MC+KITTI with and without scale transform achieve similar numerical results on test split of KITTI. This is reasonable since the KITTI data set is collected with stereo cameras that have a fixed baseline between them. On the other hand, the MC data set consists of diverse scenes with no fixed baseline between viewpoint images in each scene. For other test sets, models trained with the proposed loss function significantly outperform the models that are trained without the scale transform. Moreover, our proposed method also allows training on a combination of different data sets to generalise across scene depth ranges for indoor and outdoor scenes.



**Fig. 7.** Qualitative results on Mannequin Challenge data set.

**Table 5.** Results on two different test sets with and without proposed scale transform method. Lower is better for all error measures.

| Test set | Error Measure | Training Set | | | | Test set | Error Measure | Training Set | | | |
| | | MC | | MC + KITTI | | | | MC | | MC + KITTI | |
| | | w/o ST | ST | w/o ST | ST | | | w/o ST | ST | w/o ST | ST |
| KITTI | Abs Rel | 0.426 | 0.276 | 0.116 | **0.108** | NYU | Abs Rel | 0.282 | 0.201 | 0.333 | **0.193** |
| | Sq Rel | 4.933 | 2.563 | 0.864 | **0.855** | | RMS | 0.921 | 0.718 | 1.347 | **0.686** |
| Make3D | Abs Rel | 0.424 | 0.346 | 0.347 | **0.289** | TUM Dynamic | Abs Rel | 0.252 | 0.207 | 0.273 | **0.201** |
| | RMS | 8.771 | 7.70 | 7.64 | **6.92** | | RMS | 1.050 | **0.973** | 1.133 | 1.025 |

# 5   Conclusion and Future Work

We present a generalised self-supervised monocular depth estimation method (RealMonoDepth) that overcomes the limitation of existing self-supervised methods ([20,63,18,61,16]) that are limited to scenes with fixed scale and depth range. These methods cannot be trained on moving camera data sets due to the assumption of a fixed baseline and are unable to generalise to unseen data sets with different depth ranges. RealMonoDepth addresses all of these limitations by allowing simultaneous training on a combination of indoor and outdoor scenes with varying depth ranges. This leads to significantly improved generalisation performance across indoor and outdoor scenes and scenes which are unseen during training, and removes the dependance on a fixed camera baseline. The proposed method allows mixing stereo and moving camera data sets (MC + KITTI) improving on state-of-the-art performance in single view depth estimation across five benchmark data sets including data sets with varying depth range.

Success of deep networks depends on the use of large data sets. The proposed self-supervised training from sequences captured from a single camera allows us to train the network on diverse uncontrolled in the wild data sets, such as the Mannequin Challenge (MC) data set used in this work. This opens the door to further generalisation through training across even larger and more diverse scenes. A limitation of the proposed method is that it works with static scenes during training as with all the other single-image depth estimation methods. At test time we estimate depth from a single image and are therefore able to handle dynamic scenes with a wide variety of scale (see supplementary video). An interesting potential future work might be to extend the training of our method for dynamic scenes to increase the diversity of the data.

# References

1. Aanæs, H., Jensen, R.R., Vogiatzis, G., Tola, E., Dahl, A.B.: Large-scale data for multiple-view stereopsis. International Journal of Computer Vision pp. 1–16 (2016)
2. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al.: Tensorflow: A system for large-scale machine learning. In: 12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16). pp. 265–283 (2016)
3. Bian, J., Li, Z., Wang, N., Zhan, H., Shen, C., Cheng, M.M., Reid, I.: Unsupervised scale-consistent depth and ego-motion learning from monocular video. In: Advances in Neural Information Processing Systems. pp. 35–45 (2019)
4. Casser, V., Pirk, S., Mahjourian, R., Angelova, A.: Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 8001–8008 (2019)
5. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012 (2015)
6. Chen, P.Y., Liu, A.H., Liu, Y.C., Wang, Y.C.F.: Towards scene understanding: Unsupervised monocular depth estimation with semantic-aware representation. In:

Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2624–2632 (2019)

7. Chen, W., Fu, Z., Yang, D., Deng, J.: Single-image depth perception in the wild. In: Advances in neural information processing systems. pp. 730–738 (2016)

8. Choi, S., Zhou, Q.Y., Miller, S., Koltun, V.: A large dataset of object scans. arXiv:1602.02481 (2016)

9. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3213–3223 (2016)

10. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5828–5839 (2017)

11. Davison, A.J.: Real-time simultaneous localisation and mapping with a single camera. In: null. p. 1403. IEEE (2003)

12. Dosovitskiy, A., Fischer, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., Van Der Smagt, P., Cremers, D., Brox, T.: Flownet: Learning optical flow with convolutional networks. In: Proceedings of the IEEE international conference on computer vision. pp. 2758–2766 (2015)

13. Eigen, D., Fergus, R.: Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In: Proceedings of the IEEE international conference on computer vision. pp. 2650–2658 (2015)

14. Eigen, D., Puhrsch, C., Fergus, R.: Depth map prediction from a single image using a multi-scale deep network. In: Advances in neural information processing systems. pp. 2366–2374 (2014)

15. Fu, H., Gong, M., Wang, C., Batmanghelich, K., Tao, D.: Deep ordinal regression network for monocular depth estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2002–2011 (2018)

16. Garg, R., BG, V.K., Carneiro, G., Reid, I.: Unsupervised cnn for single view depth estimation: Geometry to the rescue. In: European Conference on Computer Vision. pp. 740–756. Springer (2016)

17. 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 (2012)

18. Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 270–279 (2017)

19. Godard, C., Mac Aodha, O., Firman, M., Brostow, G.J.: Digging into self-supervised monocular depth estimation. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3828–3838 (2019)

20. Gordon, A., Li, H., Jonschkowski, R., Angelova, A.: Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. arXiv preprint arXiv:1904.04998 (2019)

21. Guo, X., Li, H., Yi, S., Ren, J., Wang, X.: Learning monocular depth by distilling cross-domain stereo networks. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 484–500 (2018)

22. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)

23. Huang, P.H., Matzen, K., Kopf, J., Ahuja, N., Huang, J.B.: Deepmvs: Learning multi-view stereopsis. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2821–2830 (2018)

24. Jaderberg, M., Simonyan, K., Zisserman, A., et al.: Spatial transformer networks. In: Advances in neural information processing systems. pp. 2017–2025 (2015)

25. Jafari, O.H., Groth, O., Kirillov, A., Yang, M.Y., Rother, C.: Analyzing modular cnn architectures for joint depth prediction and semantic segmentation. In: 2017 IEEE International Conference on Robotics and Automation (ICRA). pp. 4620–4627. IEEE (2017)

26. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)

27. Kuznietsov, Y., Stuckler, J., Leibe, B.: Semi-supervised deep learning for monocular depth map prediction. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 6647–6655 (2017)

28. Laina, I., Rupprecht, C., Belagiannis, V., Tombari, F., Navab, N.: Deeper depth prediction with fully convolutional residual networks. In: 2016 Fourth international conference on 3D vision (3DV). pp. 239–248. IEEE (2016)

29. Li, B., Dai, Y., He, M.: Monocular depth estimation with hierarchical fusion of dilated cnns and soft-weighted-sum inference. Pattern Recognition **83**, 328–339 (2018)

30. Li, B., Shen, C., Dai, Y., Van Den Hengel, A., He, M.: Depth and surface normal estimation from monocular images using regression on deep features and hierarchical crfs. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1119–1127 (2015)

31. Li, R., Wang, S., Long, Z., Gu, D.: Undeepvo: Monocular visual odometry through unsupervised deep learning. In: 2018 IEEE International Conference on Robotics and Automation (ICRA). pp. 7286–7291. IEEE (2018)

32. Li, Z., Dekel, T., Cole, F., Tucker, R., Snavely, N., Liu, C., Freeman, W.T.: Learning the depths of moving people by watching frozen people. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4521–4530 (2019)

33. Li, Z., Snavely, N.: Megadepth: Learning single-view depth prediction from internet photos. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2041–2050 (2018)

34. Liu, F., Shen, C., Lin, G.: Deep convolutional neural fields for depth estimation from a single image. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5162–5170 (2015)

35. Liu, M., Salzmann, M., He, X.: Discrete-continuous depth estimation from a single image. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 716–723 (2014)

36. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015)

37. Mahjourian, R., Wicke, M., Angelova, A.: Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5667–5675 (2018)

38. 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: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4040–4048 (2016)

39. Montemerlo, M., Thrun, S., Koller, D., Wegbreit, B., et al.: Fastslam: A factored solution to the simultaneous localization and mapping problem. Aaai/iaai **593598** (2002)
40. Pillai, S., Ambruş, R., Gaidon, A.: Superdepth: Self-supervised, super-resolved monocular depth estimation. In: 2019 International Conference on Robotics and Automation (ICRA). pp. 9250–9256. IEEE (2019)
41. Planche, B., Wu, Z., Ma, K., Sun, S., Kluckner, S., Lehmann, O., Chen, T., Hutter, A., Zakharov, S., Kosch, H., et al.: Depthsynth: Real-time realistic synthetic data generation from cad models for 2.5 d recognition. In: 2017 International Conference on 3D Vision (3DV). pp. 1–10. IEEE (2017)
42. Poggi, M., Tosi, F., Mattoccia, S.: Learning monocular depth estimation with unsupervised trinocular assumptions. In: 2018 International Conference on 3D Vision (3DV). pp. 324–333. IEEE (2018)
43. Ricci, E., Ouyang, W., Wang, X., Sebe, N., et al.: Monocular depth estimation using multi-scale continuous crfs as sequential deep networks. IEEE transactions on pattern analysis and machine intelligence **41**(6), 1426–1440 (2018)
44. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. pp. 234–241. Springer (2015)
45. Roy, A., Todorovic, S.: Monocular depth estimation using neural regression forest. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5506–5514 (2016)
46. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. International journal of computer vision **115**(3), 211–252 (2015)
47. Saxena, A., Sun, M., Ng, A.Y.: Make3d: Learning 3d scene structure from a single still image. IEEE transactions on pattern analysis and machine intelligence **31**(5), 824–840 (2008)
48. Schönberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
49. Shilane, P., Min, P., Kazhdan, M., Funkhouser, T.: The princeton shape benchmark. In: Proceedings Shape Modeling Applications, 2004. pp. 167–178. IEEE (2004)
50. Silberman, N., Hoiem, D., Kohli, P., Fergus, R.: Indoor segmentation and support inference from rgbd images. In: European conference on computer vision. pp. 746–760. Springer (2012)
51. Steinbrücker, F., Sturm, J., Cremers, D.: Real-time visual odometry from dense rgb-d images. In: 2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops). pp. 719–722. IEEE (2011)
52. Sturm, J., Engelhard, N., Endres, F., Burgard, W., Cremers, D.: A benchmark for the evaluation of rgb-d slam systems. In: 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems. pp. 573–580. IEEE (2012)
53. Ummenhofer, B., Zhou, H., Uhrig, J., Mayer, N., Ilg, E., Dosovitskiy, A., Brox, T.: Demon: Depth and motion network for learning monocular stereo. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5038–5047 (2017)
54. Wang, C., Miguel Buenaposada, J., Zhu, R., Lucey, S.: Learning depth from monocular videos using direct methods. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2022–2030 (2018)

55. Wang, P., Shen, X., Lin, Z., Cohen, S., Price, B., Yuille, A.L.: Towards unified depth and semantic prediction from a single image. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2800–2809 (2015)

56. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P., et al.: Image quality assessment: from error visibility to structural similarity. IEEE transactions on image processing **13**(4), 600–612 (2004)

57. Wong, A., Soatto, S.: Bilateral cyclic constraint and adaptive regularization for unsupervised monocular depth prediction. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5644–5653 (2019)

58. Xiang, Y., Mottaghi, R., Savarese, S.: Beyond pascal: A benchmark for 3d object detection in the wild. In: IEEE winter conference on applications of computer vision. pp. 75–82. IEEE (2014)

59. Xu, D., Ricci, E., Ouyang, W., Wang, X., Sebe, N.: Multi-scale continuous crfs as sequential deep networks for monocular depth estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5354–5362 (2017)

60. Yao, Y., Luo, Z., Li, S., Fang, T., Quan, L.: Mvsnet: Depth inference for unstructured multi-view stereo. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 767–783 (2018)

61. Yin, Z., Shi, J.: Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1983–1992 (2018)

62. Zhan, H., Garg, R., Saroj Weerasekera, C., Li, K., Agarwal, H., Reid, I.: Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 340–349 (2018)

63. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1851–1858 (2017)

# Supplementary Material

## 1 Network Implementation Details

We use the standard pretrained resnet50 encoder, $resnet\_v1\_50.ckpt$, officially provided by Tensorflow [2]. The decoder weights are initialised randomly and the details of our architecture are shown in Table 1. Our decoder uses skip connections from the encoder [36] ($econv4, econv3, econv2, econv1$), $econv5$ is the final encoder output and estimates multi-resolution depth maps ($depth1$, $depth2$, $depth3$, $depth4$) in order to exploit both local and global information to resolve higher resolution details. Code will be released upon publication.

**Table 1. Decoder architecture of Relative DepthNet.** $Upsample()$ represents $2\times$ nearest-neighbor resizing operation and $+$ denotes channel-wise concatenation.

| Layer | Output Size | Kernel Size | Stride | Input | BatchNorm | Activation |
|---|---|---|---|---|---|---|
| upconv5 | $16 \times 32 \times 256$ | 3 | 1 | Upsample(econv5) | Yes | ReLU |
| iconv5 | $16 \times 32 \times 256$ | 3 | 1 | upconv5 + econv4 | Yes | ReLU |
| upconv4 | $32 \times 64 \times 128$ | 3 | 1 | Upsample(iconv5) | Yes | ReLU |
| iconv4 | $32 \times 64 \times 128$ | 3 | 1 | Upsample(upconv4) + econv3 | Yes | ReLU |
| depth4 | $32 \times 64 \times 1$ | 3 | 1 | iconv4 | No | Identity |
| upconv3 | $64 \times 128 \times 64$ | 3 | 1 | Upsample(iconv4) | Yes | ReLU |
| iconv3 | $64 \times 128 \times 64$ | 3 | 1 | upconv3 + econv2 + Upsample(depth4) | Yes | ReLU |
| depth3 | $64 \times 128 \times 1$ | 3 | 1 | iconv3 | No | Identity |
| upconv2 | $128 \times 256 \times 32$ | 3 | 1 | Upsample(iconv3) | Yes | ReLU |
| iconv2 | $128 \times 256 \times 32$ | 3 | 1 | upconv2 + econv1 + Upsample(depth3) | Yes | ReLU |
| depth2 | $128 \times 256 \times 1$ | 3 | 1 | iconv2 | No | Identity |
| upconv1 | $256 \times 512 \times 16$ | 3 | 1 | Upsample(iconv2) | Yes | ReLU |
| iconv1 | $256 \times 512 \times 16$ | 3 | 1 | upconv1 + Upsample(depth2) | Yes | ReLU |
| depth1 | $256 \times 512 \times 1$ | 3 | 1 | iconv1 | No | Identity |

## 2 Additional Results

We provide additional qualitative results in order to showcase the generalization ability of the proposed model trained on MC+KITTI using our novel loss function.

**Supplementary video.** We generate depth predictions with our model and compare against current state-of-the-art Monodepth2 [19] on sample YouTube videos which consist of diverse scenes with dynamic objects and varying depth range. Results are presented in the supplementary video. Note: the video scenes are monocular and were not seen during training. These videos are recorded with standard handheld monocular cameras and do not have ground-truth depth estimates. Each frame was processed independently i.e. the temporal relation is not used.

**Diverse scene images.** We also show qualitative results on the test set of DIW [7] dataset in Fig. 1, 2 and 3. These images constitute diversely rich content including indoor, natural and street scenes consisting of various objects

taken from arbitrary camera angles with uncontrolled lighting conditions and scene appearance. Results demonstrate plausible depth estimation for general scenes with performance comparable to human or previous depth estimation using supervised learning [7].
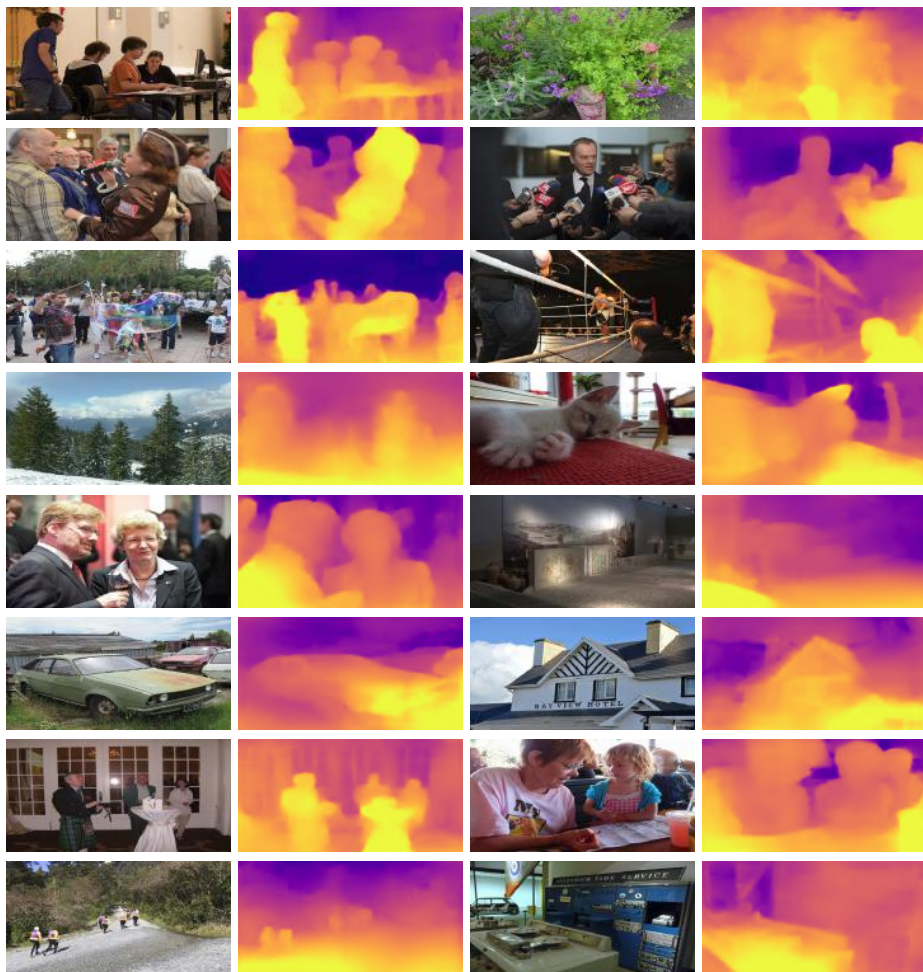


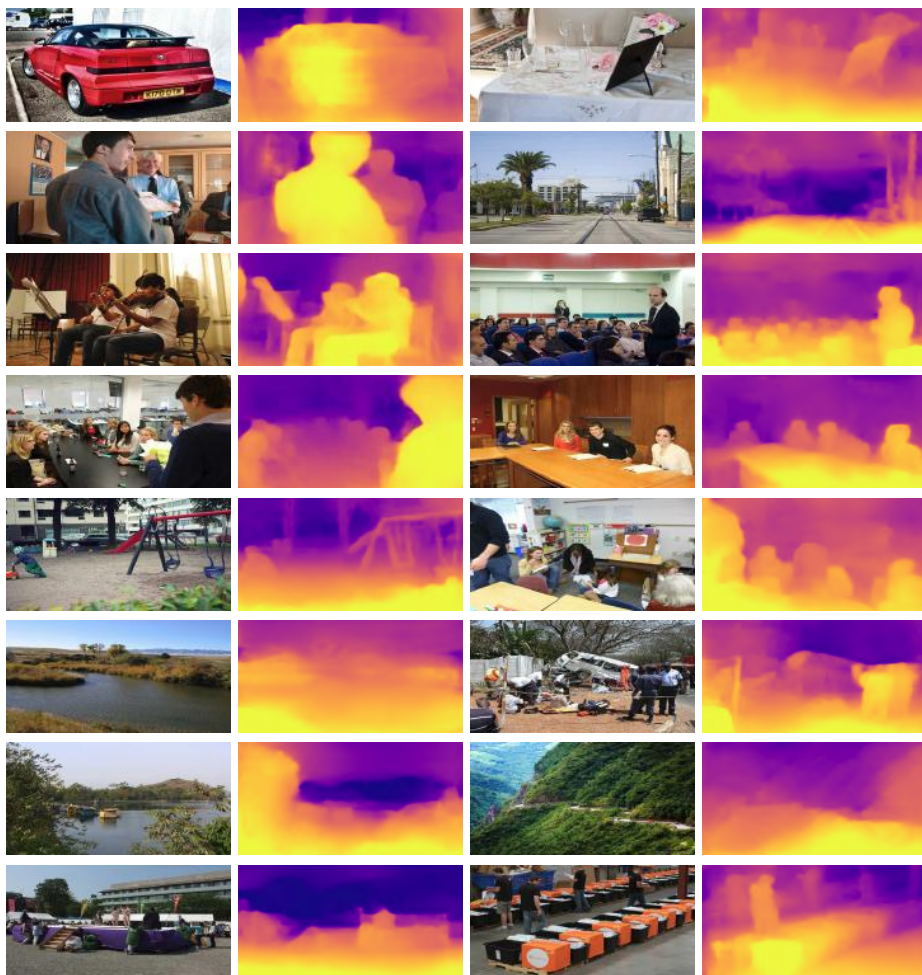**Fig. 1.** Qualitative results on the DIW test set. *(cont.)*

**Fig. 2.** Qualitative results on the DIW test set.

**Fig. 3.** Qualitative results on the DIW test set.