

Unsupervised Adversarial Depth Estimation using Cycled Generative Networks

Andrea Pilzer^{†*} Dan Xu^{‡*} Mihai Marian Puscas^{†*} Elisa Ricci^{†◊} Nicu Sebe[†]

[†]DISI, University of Trento, via Sommarive 14, Povo TN, Italy

[‡]Department of Engineering Science, University of Oxford, 17 Parks Road, Oxford, UK

[◊]Technologies of Vision, Fondazione Bruno Kessler, via Sommarive 18, Povo TN, Italy

{andrea.pilzer, mihaimarian.puscas, e.ricci, niculae.sebe}@unitn.it, danxu@robots.ox.ac.uk

Abstract

While recent deep monocular depth estimation approaches based on supervised regression have achieved remarkable performance, costly ground truth annotations are required during training. To cope with this issue, in this paper we present a novel unsupervised deep learning approach for predicting depth maps and show that the depth estimation task can be effectively tackled within an adversarial learning framework. Specifically, we propose a deep generative network that learns to predict the correspondence field (*i.e.* the disparity map) between two image views in a calibrated stereo camera setting. The proposed architecture consists of two generative sub-networks jointly trained with adversarial learning for reconstructing the disparity map and organized in a cycle such as to provide mutual constraints and supervision to each other. Extensive experiments on the publicly available datasets KITTI and Cityscapes demonstrate the effectiveness of the proposed model and competitive results with state of the art methods. The code and trained model are available: <https://github.com/andrea-pilzer/unsup-stereo-depthGAN>

1. Introduction

As one of the fundamental problems in computer vision, depth estimation has received a substantial interest in the past, also motivated by its importance in various application scenarios, such as robotics navigation, 3D reconstruction, virtual reality and autonomous driving. Over the last few years the performances of depth estimation methods have been significantly improved thanks to advanced deep learning techniques.

Most previous works considering deep architectures for predicting depth maps operate in a supervised learning setting [3, 11, 13, 24] and, specifically, devise powerful deep regression models with Convolutional Neural Networks (CNN). These models are used for monocular depth

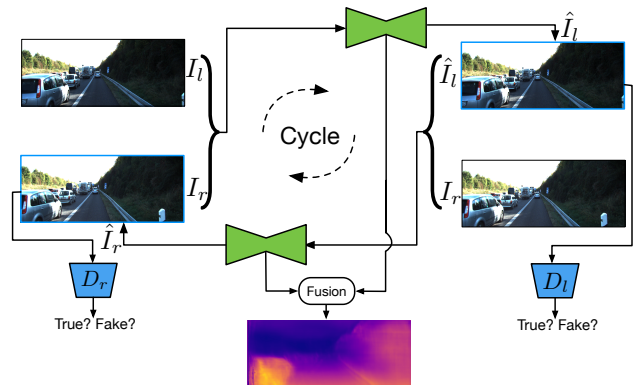


Figure 1. Motivation of the proposed unsupervised depth estimation approach using cycled generative networks optimized with adversarial learning. The left and right image synthesis in a cycle provides each other strong constraint and supervision to better optimize both generators. The \hat{I}_r and \hat{I}_l are synthesized images. Final depth estimation is obtained by fusing the output from both generators.

estimation, *i.e.* they are trained to learn the transformation from the RGB image domain to the depth domain in a pixel-to-pixel fashion. In this context, multi-scale CNN models have shown to be especially effective for estimating depth maps [3]. Upon these, probabilistic graphical models, such as Conditional Random Fields (CRFs), implemented as neural networks for end-to-end optimization, have proved to be beneficial, boosting the performance of deep regression models [13, 24]. However, supervised learning models require ground-truth depth data which are usually costly to acquire. This problem is especially relevant with deep learning architectures, as large amount of data are typically required to produce satisfactory performance. Furthermore, supervised monocular depth estimation can be regarded as an ill-posed problem due to the scale ambiguity issue [18].

To tackle these problems, recently unsupervised learning-based approaches for depth estimation have been introduced [14, 16]. These methods operate by learning the correspondence field (*i.e.* the disparity map) between the two different image views of a calibrated stereo camera us-

*The authors contributed equally in this work.

ing only the rectified left and right images. Then, given several camera parameters, the depth maps can be calculated using the predicted disparity maps. Significant progresses have been made along this research line [4, 6, 20]. In particular, Godard *et al.* [6] proposed to estimate both the direct and the reverse disparity maps using a single generative network and utilized the consistency between left and right disparity maps to constrain on the model learning. Other works proposed to facilitate the depth estimation by jointly learning the camera pose [29, 15]. These works optimized their models relying on the supervision from the image synthesis of an expected view, whose quality plays a direct influence on the performance of the estimated disparity map. However, all of these works only considered a reconstruction loss and none of them have explored using adversarial learning to improve the generation of the synthesized images.

In this paper, we follow the unsupervised learning setting and propose a novel end-to-end trainable deep network model for adversarial learning-based depth estimation given stereo image pairs. The proposed approach consists of two generative sub-networks which predict the disparity map from the left to the right view and viceversa. The two sub-networks are organized in a cycle (Fig. 1), such as to perform the image synthesis of different views in a closed loop. This new network design provides strong constraint and supervision for each image view, facilitating the optimization of both generators from the two sub-networks which are jointly learned with an adversarial learning strategy. The final disparity map is produced by combining the output from the two generators.

In summary, the main contributions of this paper are threefolds:

- To the best of our knowledge, we are the first to explore using adversarial learning to facilitate the image synthesis of different views in a unified deep network for improving the unsupervised depth estimation;
- We present a new cycled generative network structure for unsupervised depth estimation which can learn both the forward and the reverse disparity maps, and can synthesize the different image views in a closed loop. Compared with the existing generative network structures, the proposed cycled generative network is able to enforce stronger constraints from each image view and better optimize the network generators.
- Extensive experiments on two large publicly available datasets (*i.e.* KITTI and Cityscapes) demonstrate the effectiveness of both the adversarial image synthesis and the cycled generative network structure.

2. Related Work

Supervised Depth Estimation. Supervised deep learning greatly improved the performance of depth estimation.

Given enough ground-truth depth training data, deep neural networks based approaches have achieved very promising performances in recent years. Multiple large-scale depth-contained datasets [17, 19, 5, 2] have been published. In a single view setting, NYUD [17] presents indoor images while Make3D [19] is recorded in outdoors. Instead KITTI [5] and Cityscapes [2] are collected in outdoors with calibrated stereo cameras. Based on these datasets, a significant effort has been made for the supervised monocular depth estimation task [3, 13, 31, 12, 24]. The multi-scale CNN [3] and probabilistic graphical models based deep networks [13, 24, 23] also show an obvious performance boosting on the task. Xu *et al.* [25] first introduce a structured attention mechanism for learning better multi-scale deep representations for the task. However, the supervised-based approaches rely on the expensive ground-truth depth data during training, which are not flexible to deploy crossing application scenarios.

Unsupervised Depth Estimation. A more recent trend is unsupervised-based depth estimation [10, 15, 20, 28]. A remarkable advantage of unsupervised estimation lies in avoiding the use of costly ground truth depth annotations in training. Deep stereo matching models [14, 16] are proposed for direct disparity estimation. In an indirect means, Garg *et al.* [4] propose a classic approach for unsupervised monocular depth estimation based on image synthesis. Godard *et al.* [6] propose to use forward and backward reconstructions of the different image views, and multiple optimization losses are considered in the model. Zhou *et al.* [29] jointly learn the depth and the camera pose as a reinforcement in a single deep network. There are also works jointly learning the scene depth and ego-motion in monocular videos without using groundtruth data [21, 26]. However, none of these works considers the adversarial learning scheme in their models to improve the image generation quality for better depth estimation.

GANs. Generative-adversarial networks (GANs) have attracted a lot of attention for its advantage in generation problems. Godfellow *et al.* [7] revisit the generative adversarial learning strategy and show interesting results in the image generation task. After that, GANs are applied into various generation applications, and different GAN models are developed, such as CycleGAN [30] and DualGAN [27]. There are few works in the literature considering GAN models for the more challenging depth estimation task. Although Kundu *et al.* [9] investigate adversarial learning for the task, they utilize it in a context of domain adaptation in a single-track network, using a semi-supervised setting with an extra synthetic dataset, while ours considers a fully unsupervised setting and the adversarial learning in a cycled generative network aims to help the reconstruction of better image views. Both the intuition and the network design are significantly different.

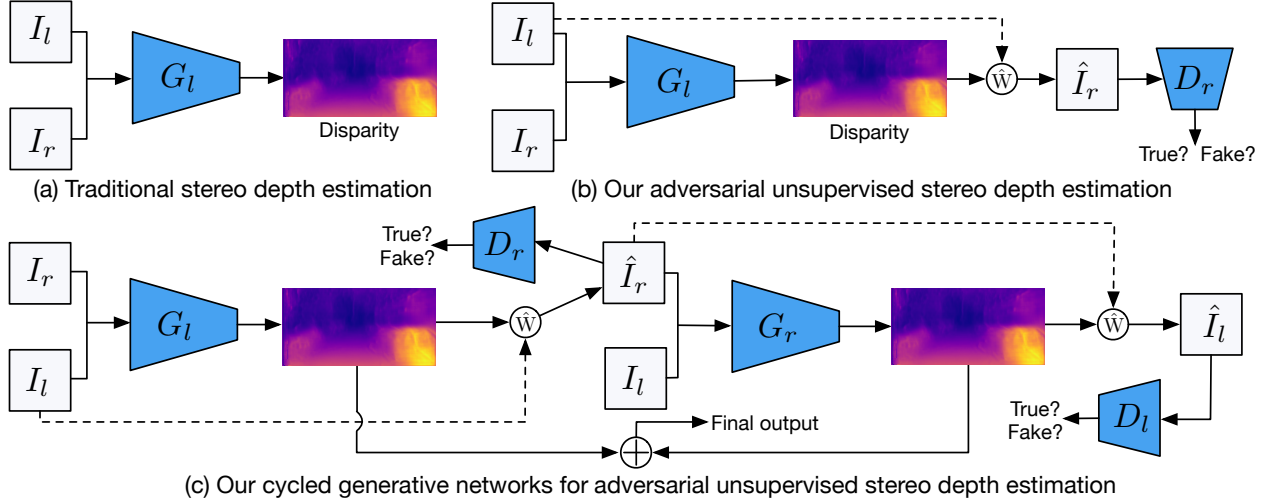


Figure 2. An illustrative comparison of different methods for unsupervised stereo depth estimation: (a) traditional stereo-matching-based depth estimation, (b) the proposed unsupervised adversarial depth estimation and (c) the proposed cycled generative networks for unsupervised adversarial depth estimation. The symbols D_l , D_r denote discriminators, and G_l , G_r denote generators. The symbol \hat{W} denotes a warping operation.

3. The Proposed Approach

We propose a novel approach for unsupervised adversarial depth estimation using cycled generative networks. An illustrative comparison of different unsupervised depth estimation models is shown in Fig. 2. Fig. 2a shows traditional stereo matching based depth estimation approaches, which basically learn a stereo matching network for directly predicting the disparity [14]. Different from the traditional stereo approaches, we estimate the disparity in an indirect means through image synthesis from different views with the adversarial learning strategy as shown in Fig. 2b. Fig. 2c shows our full model using the proposed cycled generative networks for the task. In this section we first give the problem statement, and then present the proposed adversarial learning-based unsupervised stereo depth estimation, and finally we illustrate the proposed full model and introduce the overall end-to-end optimization objective and the testing process.

3.1. Problem Statement

We target at estimating a disparity map given a pair of images from a calibrated stereo camera. The problem can be formally defined as follows: given a left image \mathbf{I}_l and a right image \mathbf{I}_r from the camera, we are interested in predicting a disparity map \mathbf{d} in which each pixel value represents an offset of the corresponding pixel between the left and the right image. If given the baseline distance b_d between the left and the right camera and the camera focal length f_l , a depth map \mathbf{D} can be calculated with the formula of $\mathbf{D} = (b_d * f_l) / \mathbf{d}$. We indirectly learn the disparity through the image synthesis. Specifically, assume that a left-to-right disparity $\mathbf{d}_r^{(l)}$ is produced from a generative network G_l

with the left-view image \mathbf{I}_l as input, and then a warping function $f_w(\cdot)$ is used to perform the synthesis of the right image view by sampling from \mathbf{I}_l , i.e. $\hat{\mathbf{I}}_r = f_w(\mathbf{d}_r^{(l)}, \mathbf{I}_l)$. A reconstruction loss between $\hat{\mathbf{I}}_r$ and \mathbf{I}_r is thus utilized to provide supervision in optimizing the network G_l .

3.2. Unsupervised Adversarial Depth Estimation

We now introduce the proposed unsupervised adversarial depth estimation approach. Assuming we have a generative network G_l composed of two sub-networks, a generative sub-network $G_l^{(l)}$ with input \mathbf{I}_l and a generative sub-network $G_l^{(r)}$ with input \mathbf{I}_r . These are used to produce two distinct left-to-right disparity maps $\mathbf{d}_r^{(l)}$ and $\mathbf{d}_r^{(r)}$ respectively, i.e. $\mathbf{d}_r^{(l)} = G_l^{(l)}(\mathbf{I}_l)$ and $\mathbf{d}_r^{(r)} = G_l^{(r)}(\mathbf{I}_r)$. The sub-network $G_l^{(l)}$ and $G_l^{(r)}$ exploit the same network structure using a convolutional encoder-decoder, where the encoders aim at obtaining compact image representations and could be shared to reduce the network capacity. Since the two disparity maps are produced from different input images, and show complementary characteristics, they are fused using a linear combination implemented as concatenation and 1×1 convolution, and we obtain an enhanced disparity map \mathbf{d}'_r , which is used to synthesize a right view image $\hat{\mathbf{I}}_r$ via the warping operation, i.e. $\hat{\mathbf{I}}_r = f_w(\mathbf{d}'_r, \mathbf{I}_l)$. Then we use an $L1$ -norm reconstruction loss \mathcal{L}_{rec} for optimization as follows:

$$\mathcal{L}_{rec}^{(r)} = \|\mathbf{I}_r - f_w(\mathbf{d}'_r, \mathbf{I}_l)\|_1 \quad (1)$$

To improve the generation quality of the image $\hat{\mathbf{I}}_r$ and benefit from the advantage of adversarial learning, we propose to use adversarial learning here for a better optimization due to its demonstrated powerful ability in the image

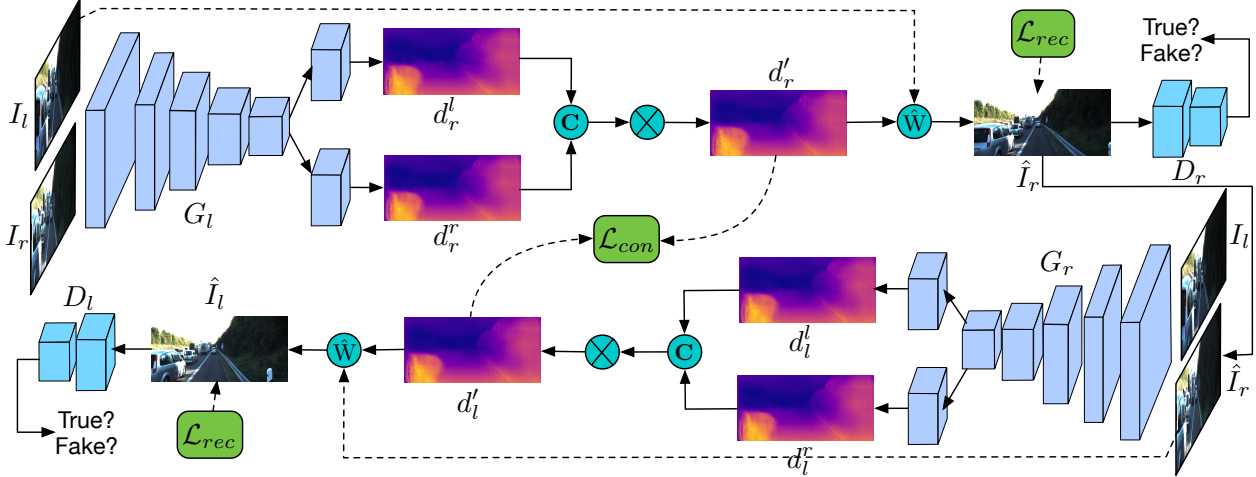


Figure 3. Illustration of the detailed framework of the proposed cyclic generative networks for unsupervised adversarial depth estimation. The symbol \odot denotes a concatenation operation; \mathcal{L}_{rec} represents the reconstruction loss for different generators; \mathcal{L}_{con} denotes a consistency loss between the disparity maps generated from the two generators.

generation task [7]. For the synthesized image $\hat{\mathbf{I}}_r$, a discriminator D_r outputting a scalar value which is used to discriminate if the image $\hat{\mathbf{I}}_r$ or \mathbf{I}_r is fake or true, and thus the adversarial objective for the generative network can be formulated as follows:

$$\mathcal{L}_{gan}^{(r)}(G_l, D_r, \mathbf{I}_l, \mathbf{I}_r) = \mathbb{E}_{\mathbf{I}_r \sim p(\mathbf{I}_r)}[\log D_r(\mathbf{I}_r)] + \mathbb{E}_{\mathbf{I}_l \sim p(\mathbf{I}_l)}[\log(1 - D_r(f_w(\mathbf{d}'_r, \mathbf{I}_l)))] \quad (2)$$

where we adopt a cross-entropy loss to measure the expectation of the image \mathbf{I}_l and \mathbf{I}_r against the distribution of the left and the right view images $p(\mathbf{I}_l)$ and $p(\mathbf{I}_r)$ respectively. Then the joint optimization loss is the combination of the reconstruction loss and the adversarial loss written as:

$$\mathcal{L}_o^{(r)} = \gamma_1 \mathcal{L}_{rec}^{(r)} + \gamma_2 \mathcal{L}_{gan}^{(r)} \quad (3)$$

where γ_1 and γ_2 are the weights for balancing the loss magnitude of the two parts to stabilize the training process. In the testing phase, the inferred \mathbf{d}'_r is the final output.

3.3. Cycled Generative Networks for Adversarial Depth Estimation

In the previous section, we presented the adversarial learning-based depth estimation approach which reconstructs from one image view to the other one in a straightforward way. In order to make the image reconstruction from different views implicitly constrain on each other, we further propose a cycled generative network structure. An overview of the proposed network structure is shown in Fig. 2. The network produces two distinct disparity maps from different view directions, and synthesizes different-view images in a closed loop. In our network design, not only the different view reconstruction loss helps for better optimization of the generators, but also the two disparity maps are connected with a consistency loss to provide strong supervision from each half cycle.

We described the half-cycle generative network with adversarial learning in Section 3.2. The cycled generative network is based on the half-cycle structure. To simplify the description, we follow the notations used in Section 3.2. Assume we have obtained a synthesized image $\hat{\mathbf{I}}_r$ from the half-cycle network, and then $\hat{\mathbf{I}}_r$ is further used as input of the next cycle generative network. Let us denote the generator as G_r , which we exploit the encoder-decoder network structure similar as G_l in Sec. 3.2. The encoder part of G_r can be also shared with the encoder of G_l to have a more compact network model (we show the performance difference between using and not using the sharing scheme), and the two distinct decoders are used to produce two right-to-left disparity maps $\mathbf{d}_l^{(l)}$ and $\mathbf{d}_l^{(r)}$ corresponding the left- and the right-view input images respectively. The two maps are also combined with the combination and the convolution operation to have a fused disparity map \mathbf{d}'_l . Then we synthesize the left-view image $\hat{\mathbf{I}}_l$ via the warping operation as $\hat{\mathbf{I}}_l = f_w(\mathbf{d}'_l, \mathbf{I}_r)$. An $L1$ -norm reconstruction loss is used for optimizing the generator G_r . Then the objective for optimizing the two generators of the full cycle writes

$$\mathcal{L}_{rec}^{(f)} = \|\mathbf{I}_r - f_w(\mathbf{d}'_r, \mathbf{I}_l)\|_1 + \|\mathbf{I}_l - f_w(\mathbf{d}'_l, \hat{\mathbf{I}}_r)\|_1 \quad (4)$$

We add a discriminator D_l for discriminating the synthesized image $\hat{\mathbf{I}}_l$, and then the adversarial learning strategy is used for both the left and the right image views in a closed loop. The adversarial objective for the full cycled model can be formulated as

$$\begin{aligned} \mathcal{L}_{gan}^{(f)}(G_l, G_r, D_r, \mathbf{I}_l, \mathbf{I}_r) = & \mathbb{E}_{\mathbf{I}_r \sim p(\mathbf{I}_r)}[\log D_r(\mathbf{I}_r)] \\ & + \mathbb{E}_{\mathbf{I}_l \sim p(\mathbf{I}_l)}[\log(1 - D_r(f_w(\mathbf{d}'_r, \mathbf{I}_l)))] + \mathbb{E}_{\mathbf{I}_l \sim p(\mathbf{I}_l)}[\log D_l(\mathbf{I}_l)] \\ & + \mathbb{E}_{\mathbf{I}_r \sim p(\mathbf{I}_r)}[\log(1 - D_l(f_w(\mathbf{d}'_l, \hat{\mathbf{I}}_r)))] \end{aligned} \quad (5)$$

Each half of the cycle network produces a disparity map corresponding to a different view translation, *i.e.* \mathbf{d}'_l and \mathbf{d}'_r . To make them constrain on each other, we add an $L1$ -norm consistence loss between these two maps as follows:

$$\mathcal{L}_{con}^{(f)} = \|\mathbf{d}'_l - f_w(\mathbf{d}'_l, \mathbf{d}'_r)\|_1 \quad (6)$$

where since the two disparity maps are for different views and are not aligned, we use the warping operation to make them pixel-to-pixel matched. The consistence loss put a strong view constraint for each half cycle and thus facilitates the learning of both half cycles.

Full objective. The full optimization objective consists of the reconstruction losses of both generators, the adversarial losses for both view synthesis and the half-cycle consistence loss. It can be written as follows:

$$\mathcal{L}_o^{(f)} = \gamma_1 \mathcal{L}_{rec}^{(f)} + \gamma_2 \mathcal{L}_{gan}^{(f)} + \gamma_3 \mathcal{L}_{con}^{(f)}. \quad (7)$$

Where $\{\gamma_i\}_{i=1}^3$ represents a set of weights for controlling the importance of different optimization parts.

Inference. When the optimization is finished, given a testing pair $\{\mathbf{I}_l, \mathbf{I}_r\}$, the testing is performed by combining the output disparity maps \mathbf{d}'_l and \mathbf{d}'_r in a weighted averaging scheme. We treat the two half cycles with equal importance, and the final disparity map \mathbf{D} is obtained as the mean of the two, *i.e.* $D = (\mathbf{d}'_l + f_w(\mathbf{d}'_l, \mathbf{d}'_r))/2$.

3.4. Network Implement Details

To describe the details of the network implementation, in terms of the generators G_l and G_r , we use a ResNet-50 backbone network for the encoder part, and the decoder part contains five deconvolution with ReLU operations in which each 2 times up-samples the feature map. The skip connections are also used to pass information from the backbone representations to the deconvolutional feature maps for obtaining more effective feature aggregation. For the discriminators D_l and D_r , we employ the same network structure which has five consecutive convolutional operations with a kernel size of 3, a stride size of 2 and a padding size of 1, and batch normalization [8] is performed after each convolutional operation. Adversarial loss is applied to output patches. For the warping operation, a bilinear sampler is used as in [6].

4. Experimental Results

We present both qualitative and quantitative results on publicly available datasets to demonstrate the performance of the proposed approach for unsupervised adversarial depth estimation.

4.1. Experimental Setup

Datasets. We carry out experiments on two large datasets, *i.e.* KITTI [5] and Cityscapes [2]. For the KITTI

dataset, we use the Eigen split [3] for training and testing. This split contains 22,600 training image pairs, and 697 test pairs. We do data augmentation with online random flipping of the images during training. The **Cityscapes** dataset is collected using a stereo camera from a driving vehicle through several German cities, during different times of the day and seasons. It presents higher resolution images and is annotated mainly for semantic segmentation. To train our model we combine the densely and coarse annotated splits to obtain 22,973 image-pairs. For testing we use the 1,525 image-pairs of the densely annotated split. The test set also has pre-computed disparity maps for the evaluation.

Parameter Setup. The proposed model is implemented using the deep learning library *TensorFlow* [1]. The input images are down-sampled to a resolution of 512×256 from 1226×370 in the case of the KITTI dataset, while for the Cityscapes dataset, at the bottom one fifth of the image is cropped following [6] and then is resized to 512×256 . The output disparity maps from two input images are fused with a learned linear combination to obtain the final disparity map with a size 512×256 . The batch size for training is set to 8 and the initial learning rate is 10^{-5} in all the experiments. We use the Adam optimizer for the optimization. The momentum parameter and the weight decay are set to 0.9 and 0.0002, respectively. The final optimization objective has weighed loss parameters $\gamma_1 = 1$, $\gamma_2 = 0.1$ and $\gamma_3 = 0.1$. The learning rate is reduced by half at both $[80k, 100k]$ steps. For our experiments we used an NVIDIA Tesla K80 with 12 GB of memory.

Detailed Training Procedure. We train the half-cycle model with a standard training procedure, *i.e.* initializing the network with random weights and making the network train for a full 50 epochs. For the cycled model we optimize the network with an iterative training procedure. After random weights initialization, we train the first half branch $\{\mathbf{I}_l, \mathbf{I}_r\} \rightarrow \hat{\mathbf{I}}_r$, with generator G_l and discriminator D_r for a 20k iteration steps. After that we train the second half branch $\{\hat{\mathbf{I}}_r, \mathbf{I}_l\} \rightarrow \hat{\mathbf{I}}_l$ with generator G_r and discriminator D_l for another 20k iterations. For the training of the first cycle branch, we do not use the cycle consistence loss since the second half branch is not trained yet. Finally we jointly train the whole network with all the losses embedded for a final round of 100k iterations.

Evaluation Metrics. To quantitatively evaluate the proposed approach, we follow several standard evaluation metrics used in previous works [3, 6, 22]. Given P the total number of pixels in the test set and \hat{d}_i , d_i the estimated depth and ground truth depth values for pixel i , we have (i) the mean relative error (abs rel): $\frac{1}{P} \sum_{i=1}^P \frac{\|\hat{d}_i - d_i\|}{d_i}$, (ii) the squared relative error (sq rel): $\frac{1}{P} \sum_{i=1}^P \frac{\|\hat{d}_i - d_i\|^2}{d_i}$, (iii) the root mean squared error (rmse): $\sqrt{\frac{1}{P} \sum_{i=1}^P (\hat{d}_i - d_i)^2}$, (iv) the mean log 10 error

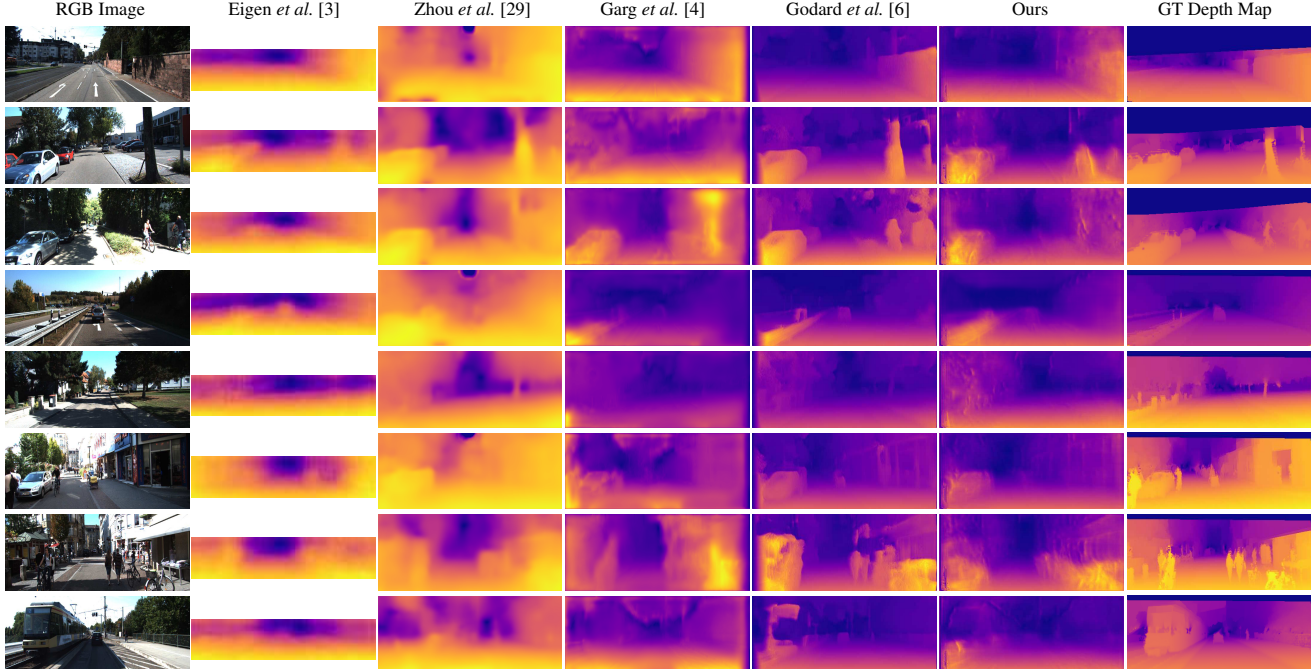


Figure 4. Qualitative comparison with different competitive approaches with both supervised and unsupervised settings on the KITTI test set. The sparse groundtruth depth maps are filled with bilinear interpolation for better visualization.

Method	Sup	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
		lower is better				higher is better		
Half-Cycle Mono	N	0.240	4.264	8.049	0.334	0.710	0.871	0.937
Half-Cycle Stereo	N	0.228	4.277	7.646	0.318	0.748	0.892	0.945
Half-Cycle + D	N	0.211	2.135	6.839	0.314	0.702	0.868	0.939
Full-Cycle + D	N	0.198	1.990	6.655	0.292	0.721	0.884	0.949
Full-Cycle + D + SE	N	0.190	2.556	6.927	0.353	0.751	0.895	0.951

Table 1. Quantitative evaluation results of different variants of the proposed approach on the KITTI dataset for the ablation study. We do not perform cropping on the depth maps for evaluation and depth range is from 0 to 80 meters.

(rmse log): $\sqrt{\frac{1}{P} \sum_{i=1}^P \|\log \hat{d}_i - \log d_i\|^2}$ (v) the accuracy with threshold t , i.e. the percentage of \hat{d}_i such that $\delta = \max(\frac{d_i}{\hat{d}_i}, \frac{\hat{d}_i}{d_i}) < t$, where $t \in [1.25, 1.25^2, 1.25^3]$.

4.2. Ablation Study

To validate the adversarial learning strategy is beneficial for the unsupervised depth estimation, and the proposed cycled generative network is effective for the task, we present an extensive ablation study on both the KITTI dataset (see Table 1) and on the Cityscape dataset (see Table 3).

Baseline Models. We have several baseline models for the ablation study, including (i) Half-cycle with a monocular setting (half-cycle mono), which uses a straight forward branch to synthesize from one image view to the other with a single disparity map output and the single RGB image is as input during testing; (ii) half-cycle with a stereo setting (half-cycle stereo), which uses a straight forward branch but

with two disparity maps produced and combined; (iii) half-cycle with a discriminator (half-cycle + D), which use a single branch as in (ii) while adds a discriminator for the image synthesis; (iv) full-cycle with two discriminators (full-cycle + D), which is our whole model using a full cycle with two discriminators added; (v) full-cycle with two discriminators and sharing encoders (full-cycle + D + SE), which has the same structure as (iv) while the parameters of the encoders of the generators are shared.

Evaluation on KITTI. As we can see from Table 1, the baseline model Half-Cycle Stereo shows significantly better performance on seven out of eight evaluation metrics than the baseline model Half-Cycle Mono, demonstrating that the utilization of the stereo images and the combination of the two estimated complementary disparity maps clearly boosts the performance.

By using the adversarial learning strategy for the image synthesis, the baseline Half-Cycle + D outperforms the

Method	Sup	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
		lower is better				higher is better		
Saxena <i>et al.</i> [18]	Y	0.280	-	8.734	-	0.601	0.820	0.926
Eigen <i>et al.</i> [3]	Y	0.190	1.515	7.156	0.270	0.692	0.899	0.967
Liu <i>et al.</i> [13]	Y	0.202	1.614	6.523	0.275	0.678	0.895	0.965
AdaDepth [9], 50m	Y	0.162	1.041	4.344	0.225	0.784	0.930	0.974
Kuznietzov <i>et al.</i> [10]	Y	-	-	4.815	0.194	0.845	0.957	0.987
Xu <i>et al.</i> [24]	Y	0.132	0.911	-	0.162	0.804	0.945	0.981
Zhou <i>et al.</i> [29]	N	0.208	1.768	6.856	0.283	0.678	0.885	0.957
Garg <i>et al.</i> [4]	N	0.169	1.08	5.104	0.273	0.740	0.904	0.962
AdaDepth [9], 50m	N	0.203	1.734	6.251	0.284	0.687	0.899	0.958
Godard <i>et al.</i> [6]	N	0.148	1.344	5.927	0.247	0.803	0.922	0.964
Ours, 80m	N	0.166	1.466	6.187	0.259	0.757	0.906	0.961
Ours with shared enc, 80m	N	0.152	1.388	6.016	0.247	0.789	0.918	0.965
Ours, 50m	N	0.158	1.108	4.764	0.245	0.771	0.915	0.966
Ours with shared enc, 50m	N	0.144	1.007	4.660	0.240	0.793	0.923	0.968

Table 2. Comparison with state of the art. Training and testing are performed on the KITTI [5] dataset. Supervised and semi-supervised methods are marked with Y in the supervision column, unsupervised methods with N. Numbers are obtained on Eigen test split with Garg image cropping. Depth predictions are capped at the common threshold of 80 meters, if capped at 50 meters we specify it.

Method	Sup	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
		lower is better				higher is better		
Half-Cycle Mono	N	0.467	7.399	5.741	0.493	0.735	0.890	0.945
Half-Cycle Stereo	N	0.462	6.097	5.740	0.377	0.708	0.873	0.937
Half-Cycle + D	N	0.438	5.713	5.745	0.400	0.711	0.877	0.940
Full-Cycle + D	N	0.440	6.036	5.443	0.398	0.730	0.887	0.944

Table 3. Quantitative evaluation results of different variants of the proposed approach on the Cityscapes dataset for the ablation study.

baseline Half-Cycle Stereo with around 1.7 points gain on the metric of Abs Rel, which verifies our initial intuition of using the adversarial learning to improve the quality of the image synthesis, and thus gain the improvement of the disparity prediction. In addition, we also observe in the training process, the adversarial learning helps to maintain a more stable convergence trend with small oscillations in terms of the training loss than the one without it (*i.e.* Half-Cycle Stereo), probably leading to a better optimized model.

It is also clear to observe that the proposed cycled generative network with adversarial learning (Full-Cycle + D) achieved much better results than the models with only half cycle (Half-Cycle + D) on all the metrics. Specifically, the Full-Cycle + D model improves the Abs Rel around 2 points, and also improves the accuracy around 1.9 points over Half-Cycle + D. The significant improvement demonstrates the effectiveness of the proposed network design, confirming that the cycled strategy brings stronger constraint and supervision to optimize the both generators. Finally, we also show that the propose cycled model using

a sharing encoder for the generator (Full-Cycle + D + SE). By using the sharing structure, we obtain even better results than the non-sharing model (Full-Cycle + D), which is probably because the shared one has a more compact network structure and thus is relatively easier to optimize with a limited number of training samples.

Evaluation on Cityscapes. We also conduct another ablation study on the Cityscapes dataset and the results are shown in Table 3. We can mostly observe similar trend of the performance gain of the different baseline models as we already analyzed on the KITTI dataset. The performance comparison of the baselines on this challenging dataset further confirms the advantage of the proposed approach. For the comparison of the model Half-Cycle + D and the model Full-Cycle + D, although the latter one achieves slightly worse results on the first two error metrics, it still produces clearly better performance on the rest six evaluation metrics. Since there is no official evaluation protocol for depth estimation on this dataset, the results are evaluated with the protocol on the KITTI, and are directly evaluated on the disparity maps as they are directly proportional to each other.

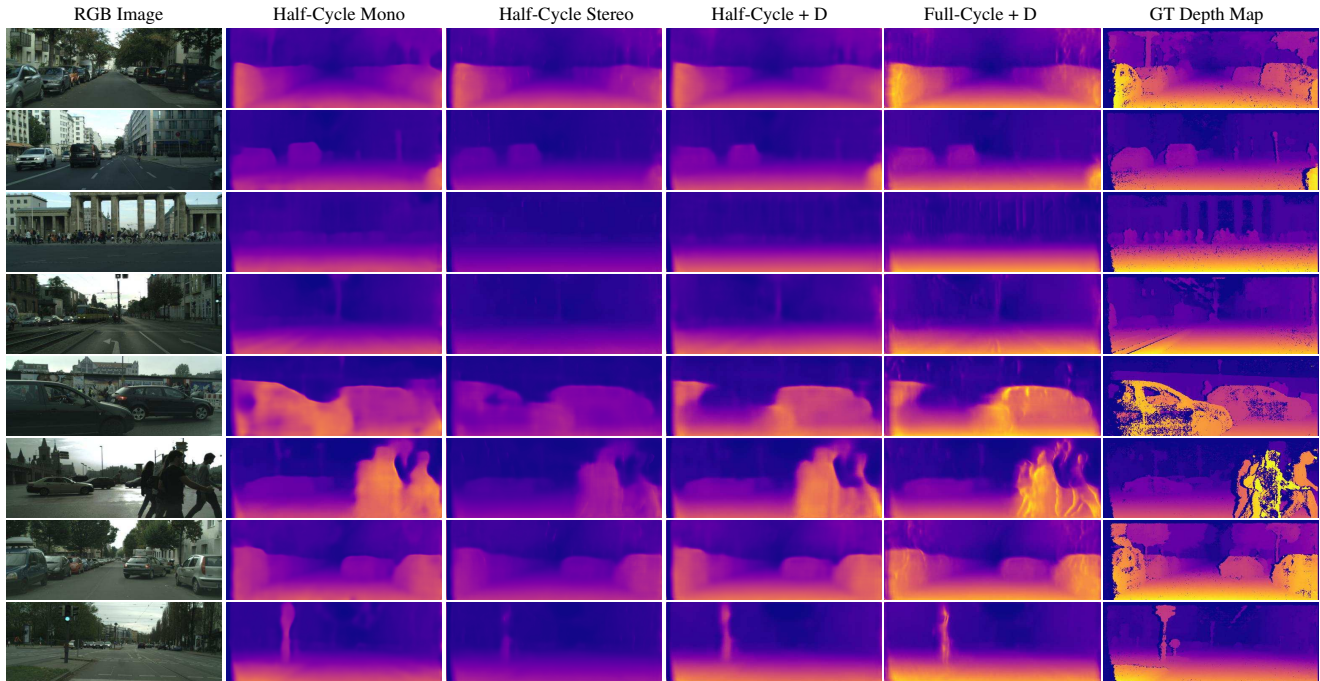


Figure 5. Qualitative comparison of different baseline models of the proposed approach on the Cityscapes testing dataset.

In Fig. 5, some qualitative comparison of the baseline models are presented.

4.3. State of the Art Comparison

In Table 2, we compare the proposed full model with several state-of-the-art methods, including the ones with the supervised setting, *i.e.* Saxena *et al.* [18], Eigen *et al.* [3], Liu *et al.* [13], AdaDepth [9], Kuznietzov *et al.* [10] and Xu *et al.* [24], and the ones with the unsupervised setting, *i.e.* Zhou *et al.* [29], AdaDepth [9], Garg *et al.* [4] and Godard *et al.* [6]. Among all the supervised approaches, we have achieved very competitive performance to the best one of them (*i.e.* Xu *et al.* [24]), while ours is totally unsupervised without using any ground-truth depth data in training. For comparison with the unsupervised methods, we are also very close to the best competitor (*i.e.* Godard *et al.* [6]). AdaDepth [9] is the most technically related to our approach, which considers adversarial learning in a context of domain adaptation with extra synthetic training data. Ours significantly outperforms their results with both the supervised and unsupervised setting, further demonstrating the effectiveness of the means we considered and proposed for unsupervised depth estimation with the adversarial learning strategy. As far as we know, there are not quantitative results presented in the existing works on the Cityscapes dataset.

4.4. Analysis on the Time Aspect.

For the training of the whole network model, on a single Tesla K80 GPU, it takes around 45 hours on KITTI dataset

with around 22k training images. For the running time, in our case with the resolution of 512×256 , the inference of one image takes around 0.140 seconds, which is a near real-time processing speed.

5. Conclusion

We have presented a novel approach for unsupervised deep learning for the depth estimation task using the adversarial learning strategy in a proposed cycled generative network structure. The new approach provides a new insight to the community that shows depth estimation can be effectively tackled via an unsupervised adversarial learning of the stereo image synthesis. More specifically, a generative deep network model is proposed to learn to predict the disparity map between two image views under a calibrated stereo camera setting. Two symmetric generative sub-networks are respectively designed to generate images from different views, and they are further merged to form a closed cycle which is able to provide strong constraint and supervision to optimize better the dual generators of the two sub-networks. Extensive experiments are conducted on two publicly available datasets (*i.e.* KITTI and Cityscapes). The results demonstrate the effectiveness of the proposed model, and show very competitive performance compared to state-of-the-arts on the KITTI dataset.

The future work would contain using attention mechanism to guide the learning of the feature representations of the generators, and also consider using the graphical models for structured prediction on the output disparity map to have predictions with better scene structures.

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [2] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.
- [3] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. In *NIPS*, 2014.
- [4] R. Garg, V. K. BG, G. Carneiro, and I. Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *ECCV*. Springer, 2016.
- [5] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012.
- [6] C. Godard, O. Mac Aodha, and G. J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, volume 2, page 7, 2017.
- [7] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [8] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [9] J. N. Kundu, P. K. Uppala, A. Pahuja, and R. V. Babu. Adadepth: Unsupervised content congruent adaptation for depth estimation. In *CVPR*, 2018.
- [10] Y. Kuznetsov, J. Stückler, and B. Leibe. Semi-supervised deep learning for monocular depth map prediction. *CVPR*, 2017.
- [11] L. Ladicky, J. Shi, and M. Pollefeys. Pulling things out of perspective. In *CVPR*, 2014.
- [12] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. In *3DV*, 2016.
- [13] F. Liu, C. Shen, G. Lin, and I. Reid. Learning depth from single monocular images using deep convolutional neural fields. *TPAMI*, 38(10):2024–2039, 2016.
- [14] W. Luo, A. G. Schwing, and R. Urtasun. Efficient deep learning for stereo matching. In *CVPR*, 2016.
- [15] R. Mahjourian, M. Wicke, and A. Angelova. Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In *CVPR*, 2018.
- [16] N. Mayer, E. Ilg, P. Hausser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *CVPR*, 2016.
- [17] P. K. Nathan Silberman, Derek Hoiem and R. Fergus. Indoor segmentation and support inference from rgbd images. In *ECCV*, 2012.
- [18] A. Saxena, S. H. Chung, and A. Y. Ng. Learning depth from single monocular images. In *NIPS*, 2006.
- [19] A. Saxena, M. Sun, and A. Y. Ng. Make3d: Learning 3d scene structure from a single still image. *TPAMI*, 31(5):824–840, 2009.
- [20] C. Wang, J. M. Buenaposada, R. Zhu, and S. Lucey. Learning depth from monocular videos using direct methods. In *CVPR*, 2018.
- [21] C. Wang, J. M. Buenaposada, R. Zhu, and S. Lucey. Learning depth from monocular videos using direct methods. In *CVPR*, 2018.
- [22] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, and A. Yuille. Towards unified depth and semantic prediction from a single image. In *CVPR*, 2015.
- [23] D. Xu, E. Ricci, W. Ouyang, X. Wang, and N. Sebe. Multi-scale continuous crfs as sequential deep networks for monocular depth estimation. In *CVPR*, 2017.
- [24] D. Xu, E. Ricci, W. Ouyang, X. Wang, and N. Sebe. Monocular depth estimation using multi-scale continuous crfs as sequential deep networks. *TPAMI*, 2018.
- [25] D. Xu, W. Wang, H. Tang, H. Liu, N. Sebe, and E. Ricci. Structured attention guided convolutional neural fields for monocular depth estimation. In *CVPR*, 2018.
- [26] Z. Yang, P. Wang, Y. Wang, W. Xu, and R. Nevatia. Every pixel counts: Unsupervised geometry learning with holistic 3d motion understanding. *arXiv preprint arXiv:1806.10556*, 2018.
- [27] Z. Yi, H. Zhang, P. Tan, and M. Gong. Dualgan: Unsupervised dual learning for image-to-image translation. *arXiv preprint*, 2017.
- [28] H. Zhan, R. Garg, C. S. Weerasekera, K. Li, H. Agarwal, and I. Reid. Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. *arXiv preprint arXiv:1803.03893*, 2018.
- [29] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, volume 2, page 7, 2017.
- [30] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017.
- [31] W. Zhuo, M. Salzmann, X. He, and M. Liu. Indoor scene structure analysis for single image depth estimation. In *CVPR*, 2015.