# Enhancing self-supervised monocular depth estimation with traditional visual odometry

Lorenzo Andraghetti[1]     Panteleimon Myriokefalitakis[1]     Pier Luigi Dovesi[1]     Belen Luque[1]
Matteo Poggi[2]     Alessandro Pieropan[1]     Stefano Mattoccia[2]

[1]Univrses AB     [2]University of Bologna

## Abstract

*Estimating depth from a single image represents an attractive alternative to more traditional approaches leveraging multiple cameras. In this field, deep learning yielded outstanding results at the cost of needing large amounts of data labeled with precise depth measurements for training. An issue softened by self-supervised approaches leveraging monocular sequences or stereo pairs in place of expensive ground truth depth annotations. This paper enables to further improve monocular depth estimation by integrating into existing self-supervised networks a geometrical prior. Specifically, we propose a sparsity-invariant autoencoder able to process the output of conventional visual odometry algorithms working in synergy with depth-from-mono networks. Experimental results on the KITTI dataset show that by exploiting the geometrical prior, our proposal: i) outperforms existing approaches in the literature and ii) couples well with both compact and complex depth-from-mono architectures, allowing for its deployment on high-end GPUs as well as on embedded devices (e.g., NVIDIA Jetson TX2).*

## 1. Introduction

Researchers have tackled the problem of estimating depth from images for decades. Understanding depth allows to interpret the environment in three dimensions and ultimately enabling the construction of 3D maps particularly useful for autonomous navigation of mobile robotics platforms or augmented and virtual reality applications.

Most of the traditional approaches to estimate depth rely on the assumption of having multiple observations of the scene, either in time (e.g structure from motion) or in space (e.g. stereo or multi-view setup), and exploit hand-crafted features to find correspondences between images to estimate sparse depth measurements of the observed scene [18].
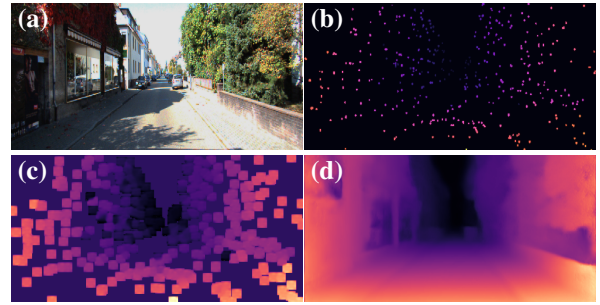


Figure 1. Monocular depth estimation enhanced by visual odometry (VO). (a) Reference image, (b) sparse 3D points by a monocular VO pipeline, (c) initial densification, (d) final depth map.

More recently, machine learning approaches have shown remarkable advances in the field [13], enabling the dense estimation of depth from a single image, given that a large amount of labelled data is available at training time.

Self-supervised paradigms relax this constraint [16, 51], replacing the need for ground truth depth annotations, usually obtained by means of active sensors [15], with additional images acquired with stereo cameras [16] or a single moving camera [51]. The former strategy is usually more effective, being both the relative pose between the two cameras and the scale factor known.

Despite the promising results achieved by depth-from-mono frameworks, they often fail in presence of ambiguous environments or elements rarely observed during the training procedure. This is caused by the absence of geometrical cues during the depth predictions, which is mostly learned upon context and semantic content of the observed scene. Inspired by the ability of humans in inferring depth from a single eye by leveraging prior knowledge (e.g. the size of known objects) [21], we propose to improve the estimation of depth by introducing a geometrical prior at inference time. Such prior comes in the form of sparse 3D measurements obtained by a traditional visual odometry (VO) algorithm that estimates the structure from consecutive images,

1

most likely scenario occurring in autonomous navigation.

Specifically, we propose a framework that combines a sparsity-invariant [41] autoencoder, which enriches our geometrical prior produced by a traditional VO algorithm, with stereo self-supervised models [16, 35] to predict depth-from-mono avoiding the need of ground truth data which is hard to obtain. Experimental results on the KITTI dataset [15] support the two main contributions of our work:

- our framework outperforms self-supervised approaches in literature.

- our strategy couples well with both complex [16] and compact [35] models, making it suited for deployment on high-end GPUs, as well as on embedded devices.

We point out that, conversely to traditional depth completion task [41] whose aim is to densify an accurate, but sparse set of depth measurements usually sourced by an active sensor such as LiDAR [41], our approach keeps the depth estimation task in the image domain and does not rely on data from any other external source.

Figure 1 shows an overview of the proposed approach: given a single image (a) and a set of 3D points obtained through VO (b), these latter are processed by the autoencoder (c) and exploited to support final depth map estimation (d).

## 2. Related Work

We briefly review the literature concerning VO, moving then to the advances in monocular depth estimation.

**Visual odometry algorithms.** Large progress has been achieved in the development of VO and SLAM methods [9, 10, 31, 32]. Although a stereo setup [9, 11, 31] avoids scale ambiguity, recent trend aims at recovering the scale of monocular VO exploiting geometry [43, 12] or deep learning [39, 47, 46].

Conversely to approaches leveraging depth to improve monocular VO and SLAM [39, 46], in this work we aim at boosting depth-from-mono accuracy exploiting VO.

**Supervised depth-from-mono.** The first approaches were supervised and they needed indeed ground truth data to enforce the network to infer depth. Among seminal works, Saxena *et al.* [37] proposed a method to estimate the absolute scales of different image patches and inferred the depth image using a Markov Random Field model, Ladick *et al.* [25] incorporated semantic segmentation into their model to improve results. With the increasing availability of ground truth depth data, supervised approaches based on CNNs [7, 27] appeared and rapidly outperformed [26, 27, 45] previous techniques. State-of-the-art in this field is DORN [13] trained with ordinal regression loss.

**Self-supervised depth-from-mono** An attractive trend concerns the possibility of learning depth-from-mono by replacing depth labels with multiple views of the sensed scene and leveraging on image synthesis to obtain supervision signals by having a loss on the reconstructed image. In general, acquiring images from a stereo camera enables a more effective training than using a single, moving camera, since the pose between frames is known. Concerning stereo supervision, Garg *et al.* [14] first followed this approach, while Godard *et al.* [16] introduced a left-right consistency loss. Other methods improved efficiency [35], deploying a pyramidal architecture, and accuracy by simulating a trinocular setup [36], including joint semantic segmentation [49] or adding adversarial term [1, 6]. In [33], a strategy was proposed to reduce further the energy efficiency of [35] leveraging fixed-point quantization. In [34] knowledge distillation and cycle consistency proved to improve results, while [40] introduces a stacked architecture, namely MonoResMatch, embodying virtual view synthesis and disparity computation and additional *proxy*-supervision self-sourced by running a traditional stereo algorithm [20]. Concerning supervision from single camera sequences, Zhou *et al.* [51] were the first to follow this direction. Their approach was improved including additional cues such as point-cloud alignment [30], differentiable Direct Visual Odometry (DVO) [42] and optical flow [52, 48]. As for stereo supervision, traditional structure-from-motion algorithms (SFM) have been used to provide additional supervision [23]. More recently, Casser *et al.* [3] introduced moving object segmentation and online refinement. Finally, few works combined the best of the two worlds, as in [50]. In particular, Yang *et al.* [46] proposed Deep Virtual Stereo Odometry (DVSO), a framework for monocular depth and ego-motion estimation trained on proxy-labels obtained from a stereo odometry algorithm. Finally, some approaches combine self-supervision with ground-truth labels from either LiDAR [24] or synthetic datasets [28, 17, 2].

As proven in prior works [40, 23, 46], we believe that traditional knowledge can provide additional cues to learning-based frameworks for monocular depth estimation. On the other hand, while existing works leverage such knowledge at training time only, we deploy a monocular VO algorithm to obtain geometrical priors to feed our network with. Being such priors sourced by a monocular setup, they are available at inference time in contrast to others available from stereo images [40, 46] and thus available at training time only.

**Depth completion.** This category covers a collection of methods with a variety of different input modalities (*e.g.*, relatively dense depth input [38] vs sparse depth measurements [29]; with color images for guidance [29] vs. without [41]). The completion problem becomes particularly challenging when the input depth image has very low density, because the inverse problem is ill-posed. One of the most popular scenario concerns with the use of 3D LiDARs, pro-
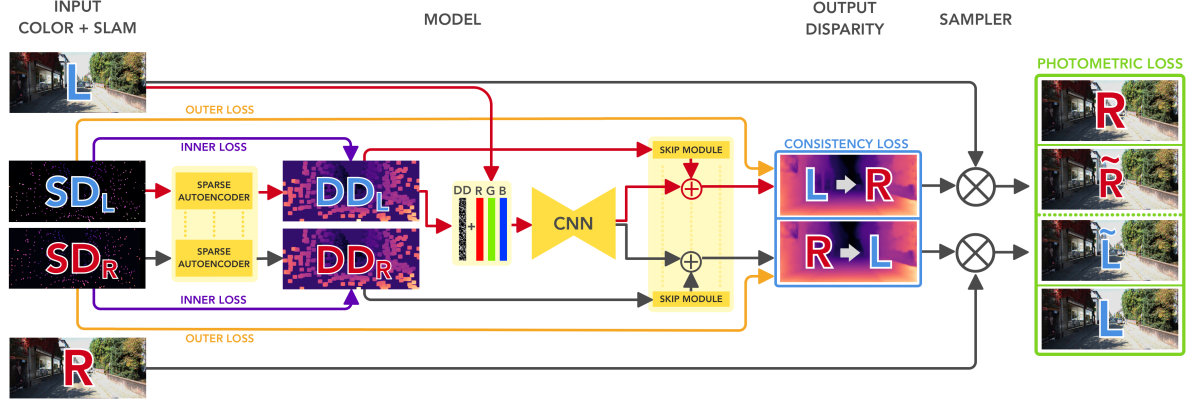
Figure 2. Overview of our framework. Sparse depths (SD) provided by a VO algorithm are fed to a sparse auto-encoder producing more dense depths (DD), forwarded then to the main network together with color image. Self-supervision is obtained by means of stereo image reprojection plus consistency (green and lightblue) and the sparse points themselves (orange and violet). At deployment, monocular cues only are required (red path)

viding roughly 5% pixels when reprojected on images [41]. Specifically, Ma and Karaman [29] proposed an end-to-end deep regression model for depth completion. Uhrig *et al.* [41] proposed sparse convolution, a variant of regular convolution operations capable of dealing with data sparsity in neural networks. Eldesokey *et al.* [8] improved the normalized convolution for confidence propagation. Chodosh *et al.* [4] incorporated the traditional dictionary learning with deep learning into a single framework for depth completion. All learning based methods [29, 41, 8, 4] are trained in a supervised manner deploying depth labels.

In contrast we leverage depth estimates obtained by means of a VO algorithm, distinguishing our approach from reviewed depth completion models usually exploiting LiDAR points that are i) sourced from very accurate, active sensors and ii) have an average density of 5% with respect to the entire image, while the VO pipeline used in our experiments only provides about 0.06% (*i.e.* 1 every 1600+ pixels).

## 3. Method

In this section, we introduce the rationale behind the proposed method and the modules deployed in our framework. We argue that it is unlikely that the entire life-cycle of an application is constrained, on most cases, to a single image acquisition at a single time frame. A popular example is represented by the autonomous driving scenario, where continuous image acquisition by means of a single camera is necessary. We aim at improving monocular depth estimation by leveraging this assumption in order to recover the geometry that is missing from a single image acquisition. For this purpose, we choose traditional VO algorithms to obtain a set of 3D points used as additional input to our framework to guide it towards more accurate estimations. In

particular, sparse 3D points are mapped to image pixels and converted to an equivalent representation with respect to the one of the final depth output. For instance, in case of stereo self-supervision [16] 3D points' depth is back-triangulated to disparity according to the specific setup (*i.e.*, baseline and focal length) deployed for training. Figure 2 shows our pipeline, made of two main modules: a sparsity-invariant autoencoder, processing the aforementioned 3D points to obtain more dense priors, and a depth estimator that outputs the final depth map when fed with the reference image and densified priors. While stereo images are required at training time, only the monocular input is processed at deployment (connected by the red path in the figure). In order to provide meaningful information to the network, the input cues are scale-aware. This can be easily obtained at training time by running a stereo VO algorithm [31], while at test time a monocular VO framework with any scale recovery, as for instance [43, 12], is required.

### 3.1. Sparsity-invariant autoencoder

The first step in our pipeline consists in processing the 3D points retrieved by means of VO. Because of their sparse nature, we design a sparsity-invariant autoencoder, since traditional convolutions results in poor performance when dealing with such kind of data, as proven in [41]. As shown in Figure 2, our autoencoder obtains a more dense depth map by means of five sparsity-invariant convolutional layers. The output of this module, namely $DD$ in the figure, is supervised by the *inner* loss shown in the figure and better described in the remainder. Figure 3 shows how the autoencoder is composed: 4 sparse-convolution layers with decreasing kernel size (9, 5, 3, 3), each one with 16 filters and stride fixed to 1 in order to keep the same resolution of the input. One final sparse-convolution pixel-wise filter is added in order to get an image that represents a denser
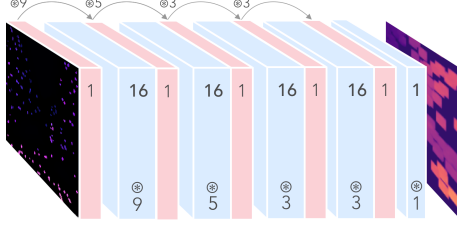
Figure 3. Structure of the sparsity-invariant autoencoder. Four convolutional layers extract 16 features each, respectively with $9 \times 9, 5 \times 5, 3 \times 3$ and $3 \times 3$ kernels.

disparity map which is used for the inner loss. Then, it is concatenated to the input image and forwarded both to the main depth estimator and to a skip residual module that will be further discussed. Since the output of the VO system is a set of sparse 3D points, it is possible to reproject them onto both left and right camera planes, generating sparse disparity maps $SD_L$ and $SD_R$. During training we employ two autoencoders with shared weights to generate both $DD_L$ and $DD_R$ from left and right sparse ones. Therefore, we enforce consistency in the losses keeping the whole system symmetric. The rationale behind this choice will be discussed shortly, while ablation experiments will highlight the contribution introduced by such strategy. This symmetry is employed during training only, while at test time a single autoencoder processes left sparse disparity map $SD_L$ to generate $DD_L$ which is given to the depth estimator after a concatenation with RGB color image.

### 3.2. Depth estimator

The recent literature provides multiple architectures for self-supervised monocular depth estimation [16, 35, 40]. To prove that our proposal is compatible with both complex and compact networks, making it suited for a wide range of applications on both high-end GPUs as well as on low-power devices, we choose two main models for our experiments: monodepth by Godard *et al.* [16] and PyD-Net by Poggi *et al.* [35]. The main difference between the two consists into the backbone used for features extraction.

The former represents the first model proposed for self-supervised monocular depth estimation from stereo images. In its more accurate version, it consists of a ResNet50 [19] encoder and about 58 million parameters. The latter deploys a compact, pyramidal features extractor, counting fewer than 2 million parameters and dramatically reducing both memory and runtime requirements [35]. Experimental results will show how the proposed pipeline is compatible with different architectures designed to maximize either accuracy [16] or efficiency [35].

In Figure 2, the aforementioned architecture (yellow block) is fed with an RGB image and densified depth cues and outputs two inverse depth maps $D^{LR}$ and $D^{RL}$, *i.e.* disparity maps, aligned respectively with the input image (left

frame of a training stereo pair) and the additional one used for supervision (right frame).

### 3.3. Skip module

In order to lighten the estimation task of the depth estimator, we add a residual skip module, further processing $DD$. This module is made of a single ResNet block [19], built by three layers respectively with kernels $1 \times 1, 3 \times 3, 1 \times 1$ and extracting 16, 16 and 64 features. In parallel, a skip connection made by a single $1 \times 1$ layer produces 64 features summed to those extracted by the latter of the previous three layers. A final $1 \times 1$ layer produces a residual correction $DD'$. For symmetry, both $DD_L, DD_R$ are processed by a shared skip module to obtain $DD'_L, DD'_R$.

Finally, we obtain two maps $d^L$ and $d^R$ as last output, respectively summing $D^{LR}$ to $DD'_L$ and $D^{RL}$ to $DD'_R$. An *outer* loss is computed between these final outputs and SD making the depth estimator, in other words, focusing on the remaining portions of the image for which no prior depth is available. Since both $d^L$ and $d^R$ are optimized, the symmetry kept by the autoencoder avoids unbalancing between losses computed on the two, as explained in Section 3.1.

### 3.4. Training Loss

Following successful attempts in literature [16, 35, 40], we deploy a multi-component loss function defined as

$$\mathcal{L} = \alpha_{st} \sum_{s=1}^{4} \mathcal{L}_{st}^s + \alpha_{in}\mathcal{L}_{in} + \alpha_{out}\mathcal{L}_{out} \tag{1}$$

where $\mathcal{L}_{st}$, $\mathcal{L}_{in}$ and $\mathcal{L}_{out}$ are respectively stereo self-supervision, inner and outer losses.

#### 3.4.1 Stereo self-supervision

We train our network using stereo self-supervision [14]. At training time, each sample consists of a stereo pair made of L and R, respectively, the image input to the model and the one used for image reprojection and the subsequent loss computation. According to Equation 1, at each scale we compute $\mathcal{L}_{st}$ as

$$\mathcal{L}_{st} = \beta_{ap}(\mathcal{L}_{ap}^L + \mathcal{L}_{ap}^R) + \beta_{ds}(\mathcal{L}_{ds}^L + \mathcal{L}_{ds}^R) \\ + \beta_{lr}(\mathcal{L}_{lr}^L + \mathcal{L}_{lr}^R) + \beta_o(\mathcal{L}_{occ}^L + \mathcal{L}_{occ}^R) \tag{2}$$

respectively made of appearance matching, disparity smoothness, left-right consistency and occlusion terms.

**Appearance Matching Loss.** enforces the reconstructed image to appear similar to the corresponding training input, combination of L1 and single scale Structured Similarity Index Measure (SSIM) [44] which compares, for each pixel of coordinates $(i, j)$, the input image $I^L$ and its reprojected

$\widetilde{I}^L$ obtained by means of bilinear warping according to disparity estimations.

$$\mathcal{L}_{ap}^L = \frac{1}{N}\sum_{ij}\gamma\frac{1-SSIM(I_{ij}^L,\widetilde{I}_{ij}^L)}{2}+(1-\gamma)||I_{ij}^L-\widetilde{I}_{ij}^L|| \tag{3}$$

where $N$ is the number of pixels and $\gamma = 0.85$.

**Disparity Smoothness Loss** enforces smooth disparities exploiting an L1 penalty on the disparity gradients $\partial d$, weighted by an edge aware term from the image.

$$\mathcal{L}_{ds}^L = \frac{1}{N}\sum_{ij}|\partial_x d_{ij}^L|e^{-||\partial_x I_{ij}^L||}+|\partial_y d_{ij}^L|e^{-||\partial_y I_{ij}^L||} \tag{4}$$

**Left-Right Consistency Loss** enforces the left and the right disparities to be consistent by using an L1 penalty between the left-to-right disparity map and the reconstructed one which comes from sampling the right-to-left in a similar manners as for the left and right images:

$$\mathcal{L}_{lr}^L = \frac{1}{N}\sum_{ij}|d_{ij}^L - d_{ij+d_{ij}^L}^R| \tag{5}$$

**Occlusion Loss** discourages artifacts near occlusions [46] by minimizing the sum of all disparities

$$\mathcal{L}_{occ}^L = \frac{1}{N}\sum_{ij}d_{ij}^L \tag{6}$$

#### 3.4.2 Inner loss

The purpose of sparsity-invariant autoencoder is to provide the depth estimator with more dense depth priors. To this aim, we enforce the output map $DD$ to be consistent with the input cues $SD$ where these are defined

$$\mathcal{L}_{in} = \frac{1}{N}\sum_{ij}|DD_{ij} - SD_{L_{ij}}| \tag{7}$$

For symmetry, this is carried out on both $DD_L$ and $DD_R$.

#### 3.4.3 Outer loss

The final prediction $d^L$ by our network is a sum of $DD'_L$ produced by the skip module and $D^{LR}$ by the depth estimator. Again, since we want to preserve the information sourced by VO, we apply a second, outer loss to enforce consistency between $SD$ and the final output

$$\mathcal{L}_{out} = \frac{1}{N}\sum_{ij}|d_{ij} - SD_{ij}| \tag{8}$$

As for the inner loss, this is carried out on both $d^L$ and $d^R$ as well for symmetry.

## 4. Experimental results

In this section, we describe the dataset and the implementation details, and report results concerning our framework in various training/testing configurations, showing that our approach is consistently beneficial to traditional self-supervised approaches. To conform to the literature, we assess the performance of monocular depth estimation techniques following the protocol by Eigen *et al.* [7], extracting data from the KITTI [15] dataset and using sparse LiDAR measurements as ground truth for evaluation.

### 4.1. Datasets

For all our experiments we compute standard metrics [7] in the field of monocular depth estimation. Abs rel, Sq rel, RMSE and RMSE log represent error measures, while $\delta < \varepsilon$ represents the percentage of predictions whose maximum between ratio and inverse ratio with respect to the ground truth is lower than $\varepsilon$, traditionally set to $1.25, 1.25^2$ and $1.25^3$.

For our experiments we use the KITTI dataset [15]. It consists of about 42382 rectified stereo pairs grouped into 61 sequences, with an image resolution of $1242 \times 375$. During acquisition, a LiDAR sensor gathered sparse depth measurements. For our experiments, we divide the entire dataset into a training and a testing set, following the traditional split by Eigen *et al.* [7]. In particular, since our method is coupled with VO algorithms, we define different subdivisions, still compliant with the Eigen split:

- for training, we adopt the $K_r, K_o$ sets introduced by Yang *et al.* [46], being the latter part of sequences 01, 02, 06, 08, 09 and 10 from KITTI odometry dataset.

- for testing, we adopt sequences 00, 04, 05 and 07 from KITTI odometry, that partially overlaps with the Eigen testing set.

This split allows for full deployment of VO algorithms both at training time, described in detail in the remainder, as well as for evaluation. Focusing on the testing split, the one we introduce counts 8691 frames, in contrast with the original one by Eigen *et al.* [7] counting only 697 images, yet being fully consistent with it (*i.e.*, there is no overlap between the 8691 frames with any of the frames from the Eigen original training set). This allows for a fair comparison with any method proposed in literature, if trained on the Eigen split and whose weights are provided by the authors, without the need for retraining.

### 4.2. Implementation details

Our framework is implemented using the Tensorflow library. We designed two variants, namely **VOMonodepth** and **VOPyD-Net**, respectively built around the depth estimators proposed by Godard *et al.* [16] and Poggi *et al.*

| | Lower is better | | | | Higher is better | | |
|---|---|---|---|---|---|---|---|
| Model | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ |
| Monodepth - ResNet [16] | 0.108 | 0.679 | 4.123 | 0.194 | 0.868 | 0.952 | 0.978 |
| baseline | 0.109 | 0.660 | 4.077 | 0.195 | 0.866 | 0.952 | 0.979 |
| + sparse-autoencoder | 0.099 | 0.666 | 3.910 | 0.200 | 0.878 | 0.947 | 0.973 |
| + sparse-autoencoder + skip | 0.099 | 0.654 | 3.843 | 0.200 | 0.879 | 0.948 | 0.973 |
| + sparse-autoencoder + skip + sym. | 0.095 | 0.621 | 3.827 | 0.186 | 0.885 | 0.952 | 0.977 |
| + sparse-autoencoder + skip + sym. + ft | **0.091** | **0.548** | **3.690** | **0.181** | **0.892** | **0.956** | **0.979** |

Table 1. Ablation experiments for VOMonodepth on KITTI [15] odometry sequences from the Eigen split [7] (8691 frames).

[35]. In both cases, at training time we fed the network with batches of 8 images, using Adam Optimizer [22] with $\beta_1 = 0.9, \beta_2 = 0.999$ and $\epsilon = 10^{-8}$ and a learning rate of $\lambda = 10^{-4}$, halved twice after $\frac{3}{5}$ and $\frac{4}{5}$ of the total epochs. The weights in our loss function have been tuned respectively to $\alpha_{st} = 1, \alpha_{in} = 5, \alpha_{out} = 2, \beta_{app} = 1, \alpha_{lr} = 1$, $\alpha_{occ} = 0.01$ and $\beta_{ds} = 0.1/r$, being $r$ the downsampling factor of each scale. According to the chosen depth estimator, we run different training schedules: for VOPyD-Net, we run 200 epochs halving the learning rate at 120 and 160, while for VOMonodepth we run 50 epochs and halve at 30 and 40, following in both cases the guidelines suggested by the authors of PyD-Net and Monodepth respectively. We perform established data augmentation procedures [16], consisting of horizontal flip of the input and color augmentation with with a 50% chance, random gamma, brightness and color shifts. At inference time, the same post-processing from [16] is applied by VOMonodepth. Images are resized to $256 \times 512$ and VO points are reprojected accordingly, except for VOPyD-Net where SD maps are provided at half the resolution for the sake of speed, then estimated DD are upsampled to the original resolution.

In our experiments, we deploy two VO odometry algorithms for training. The two respectively exploit stereo and monocular sequences. While the former provides the correct scale, the second requires a scale recovery mechanisms as for instance in [43, 12]. In order to ease the learning process, we perform a first round of training on $K_r + K_o$ feeding the network with the 3D points from stereo VO. Then, we run a further round on $K_o$ switching to the monocular VO used at testing time. We will show in the experiments how this strategy is beneficial to our framework.

For stereo VO we use ORB-SLAM2 [31]. As monocular VO algorithm with scale recovery, we use the pipeline developed by Zenuity[1] by their kind concession, the same is deployed for inference in our evaluation keeping our method fully monocular at deployment.

### 4.3. Ablation studies

We first run a set of experiments to study the impact of each design choice introduced in our framework. Pur-

posely, we train VOMonodepth in five configurations obtained by: i) directly feeding the depth estimator with VO input (*baseline*) ii) introducing the sparse-autoencoder iii) adding the skip module iv) performing symmetric training on the sparse points v) fine-tuning on $K_o$ and monocular VO.

Table 1 collects the outcome of this evaluation, carried out on the testing split mentioned above and made of 8691 frames. For comparison, we report the results achieved by Monodepth-ResNet [16]; for all models, we perform the post-processing step introduced in [16]. We point out that, while the baseline barely outperforms [16] on some metrics, the introduction of the sparse-autoencoder is crucial to boost the effectiveness of our approach. Adding the skip module (+skip) to our architecture enables for a slight improvement on Sq Rel, RMSE and $\delta$ scores. A major contribution is given by adding the symmetric training (+sym.), optimizing on VO points aligned both on the left and right images. It is worth noting that the aforementioned three configurations have been trained to leverage VO input from a stereo algorithm, while at deployment such cues comes from a monocular VO algorithm. Although the nature of input VO differs between training and testing, our technique is effective indeed at improving monocular depth estimation. Finally, by running a fine-tuning (+ft) switching from stereo to monocular VO input allows to a further, major boost on all metrics, leading to the best configuration.

### 4.4. Comparison with state-of-the-art

Table 2 reports results on the same testing split defined before. We point out once more that, since compliant with the Eigen split [7], we can compare our proposal to most existing methods. Specifically, we report in the table competitors for whose the source code or trained models are available, self-supervised either using monocular or stereo images. Unfortunately, code and models are not available for [46, 34], and hence, we are not able to compare with them. Methods marked with * have been pre-trained on CityScapes dataset [5], for whose the authors do not provide weights trained on KITTI only. The table is divided into three portions, from top to bottom: i) monocular supervised, ii) lightweight stereo-supervised and iii) complex stereo-supervised models. Our variants, VOPyD-Net and

---

[1] https://www.zenuity.com/

| Method | Supervision | PP | VO | Lower is better | | | | Higher is better | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ |
| SfmLearner [51] | Mono* | | | 0.175 | 1.309 | 5.515 | 0.247 | 0.740 | 0.916 | 0.971 |
| Vid2depth [30] | Mono* | | | 0.143 | 0.827 | 4.702 | 0.213 | 0.812 | 0.943 | 0.979 |
| GeoNet - ResNet [48] | Mono* | | | 0.141 | 0.842 | 4.688 | 0.209 | 0.812 | 0.945 | 0.981 |
| LKVO [42] | Mono* | | | 0.135 | 0.812 | 4.246 | 0.200 | 0.836 | 0.952 | 0.982 |
| DF-Net [52] | Mono* | | | 0.131 | 0.706 | 4.365 | 0.196 | 0.831 | 0.952 | 0.984 |
| Struct2depth (M) [3] | Mono* | | | 0.135 | 0.792 | 4.356 | 0.197 | 0.836 | 0.955 | 0.984 |
| PyD-Net [35] | Stereo | | | 0.130 | 0.833 | 4.569 | 0.219 | 0.825 | 0.938 | 0.974 |
| VOPyD-Net | Stereo | | ✓ | 0.105 | 0.916 | 4.916 | 0.203 | 0.874 | 0.946 | 0.974 |
| PyD-Net [35] | Stereo | ✓ | | 0.123 | 0.733 | 4.333 | 0.210 | 0.834 | 0.943 | 0.976 |
| VOPyD-Net | Stereo | ✓ | ✓ | **0.102** | **0.611** | **3.810** | **0.188** | **0.876** | **0.952** | **0.979** |
| Monodepth - ResNet [16] | Stereo | ✓ | | 0.108 | 0.679 | 4.123 | 0.194 | 0.868 | 0.952 | 0.978 |
| 3Net - ResNet [36] | Stereo | ✓ | | 0.106 | 0.627 | 3.982 | 0.192 | 0.869 | 0.953 | 0.979 |
| MonoResMatch [40] | Stereo+SGM | ✓ | | 0.102 | 0.563 | 3.725 | 0.183 | 0.885 | **0.964** | **0.986** |
| VOMonodepth - ResNet | Stereo | ✓ | ✓ | **0.091** | **0.548** | **3.690** | **0.181** | **0.892** | 0.956 | 0.979 |

Table 2. Evaluation on KITTI [15] odometry sequences from the Eigen split [7] (8691 frames). * means pre-training on CityScapes [5].

VOMonodepth, belong respectively to categories ii) and iii).

Starting from compact models, we evaluate PyD-Net [35] and its VO variant with and without applying the post-processing (PP) step introduced in [16]. However, it is worth observing that since PP requires to forward the input image twice, it adds a non-negligible overhead that is undesirable in case of deployment on embedded systems or when targeting the maximum efficiency [35]. VOPyD-Net outperforms PyD-Net by a notable margin on most metrics, except Sq Rel and RMSE. In particular, $\delta < 1.25$ receives the highest improvements, *i.e.* 87.5% pixels are below the threshold versus the 82.5 of PyD-Net. By running the post-processing, VOPyD-Net consistently outperforms it on all metrics by a significant margin. Moreover, this model also outperforms Monodepth-ResNet on all metrics and 3Net-ResNet on all except $\delta < 1.25^2$ and $\delta < 1.25^3$, despite the much lower complexity of the network and the lower runtime, as discussed further.

By coupling our strategy with a more complex architecture, as in the case of VOMonodepth-ResNet, we can outperform even MonoResMatch [40] which deploys a more accurate architecture and leverages additional supervision from SGM algorithm [20] at training time. This experiment further proves the effectiveness of our strategy, outperforming state-of-the-art methods for self-supervised monocular depth estimation.

## 4.5. Runtime analysis

We benchmark the performance of our VO variants and traditional models on two very different hardware platforms: an NVIDIA 2080 Ti GPU, having 250W of power consumption, and a Jetson TX2 module with a maximum average consumption of about 15W. The latter device represents an appealing platform for a wide range of applications. In all experiments, the TX2 board was configured for maximum performance. Table 3 collects the results of this anal-

| Method | PP | TX2 (Fps) | 2080Ti (Fps) |
|---|---|---|---|
| PyD-Net [35] | | 24.22 | 195.34 |
| VOPyD-Net | | 18.48 | 143.04 |
| Monodepth - ResNet [16] | | 3.41 | 85.12 |
| VOMonodepth - ResNet | | 2.77 | 67.13 |
| VOPyD-Net | ✓ | 8.35 | 100.11 |
| Monodepth - ResNet [16] | ✓ | 2.24 | 62.21 |
| 3Net-ResNet[36] | ✓ | 2.10 | 49.24 |
| VOMonodepth - ResNet | ✓ | 1.70 | 41.84 |
| MonoResMatch [40] | ✓ | 1.23 | 29.79 |

Table 3. Runtime comparison (averaged over 200 frames) between previous models and VO variants on Jetson TX2 and 2080Ti.

ysis, comparing Monodepth and PyD-Net with their VO counterparts. We report Fps both enabling and disabling post-processing.

Focusing on the TX2 platform, we can notice how the difference between PyD-Net and VOPyD-Net is about 30%, running respectively at 24 and 18 Fps, if post-processing is disabled, about $7\times$ and $5\times$ faster than [16] (3.41). A similar overhead, about 20%, is introduced comparing VOMonodepth-ResNet to Monodepth-ResNet. Enabling PP, VOPyD-Net still runs at more than 8 Fps, but it produces better results on each metrics (see Table 2) compared to Monodepth-ResNet, despite running more than $3\times$ faster. It also outperforms on most metrics 3Net-ResNet, running almost $4\times$ faster. VOMonodepth achieves much higher accuracy at the cost of a further drop in speed (below 2 Fps), but still 40% faster than MonoResMatch.

Switching to NVIDIA 2080Ti GPU, a similar overhead between each model and its VO variant can be noticed. Even enabling post-processing, VOPyD-Net still runs at about 100 Fps versus the 62 and 49 by Monodepth-ResNet and 3Net-ResNet, while VOMonodepth-ResNet reaches
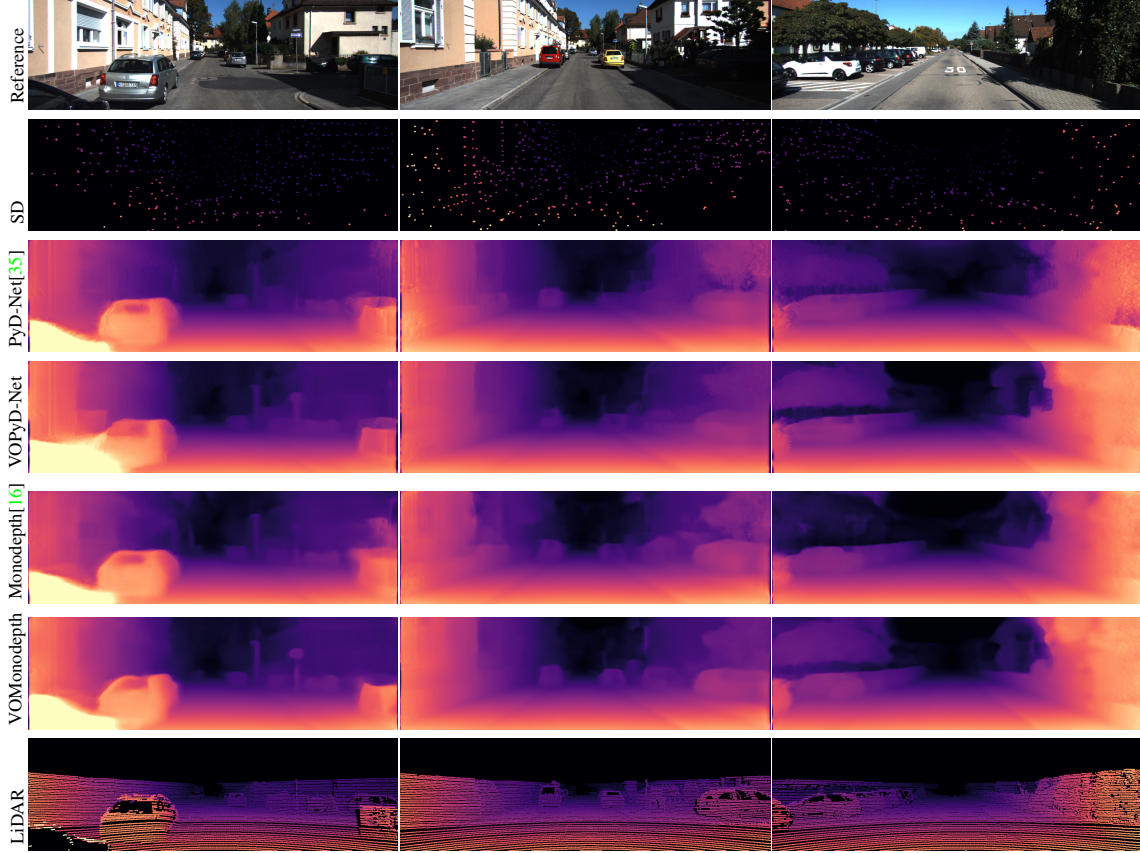
Figure 4. Qualitative results on KITTI dataset. On each column, from top to bottom: reference image, sparse VO points, depth map outputs from PyD-Net [35], VOPyD-Net, Monodepth [16], VOMonodepth and LiDAR points used for evaluation.
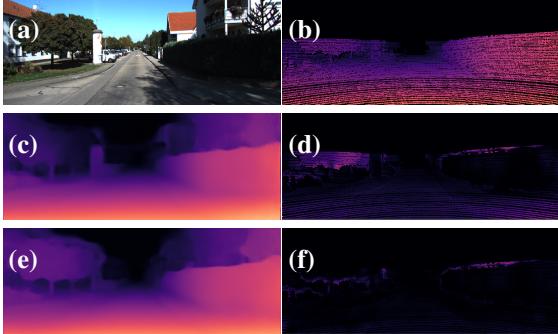


Figure 5. Qualitative comparison on KITTI dataset. (a) input image, (b) LiDAR points, depth and error maps respectively by Monodepth (c,d) and VOMonodepth (e,f).

about 40 Fps versus the 29 of MonoResMatch, making the former the best choice when high-end GPUs are available for deployment thanks to its superior accuracy.

Finally, Figure 4 shows some qualitative examples of depth map inferred by both PyD-Net and Monodepth as well as their VO counterparts. We can notice how sparse inputs improve, for instance, estimates on thin structures (left column). Figure 5 shows a comparison between Mon-odepth [16] and VOMonodepth, reporting error maps (d), (f) to highlight how the former fails at estimating the depth for trees on the left, whereas our method is successful.

## 5. Conclusion

In this paper, we have proposed a novel framework that takes into account prior knowledge to improve monocular depth estimation. We have introduced a geometrical prior obtained by estimating the movement of the camera, as it commonly happens in an autonomous navigation scenario. Our network is able to leverage on the sparse 3D measurements of a VO algorithm to improve depth accuracy. Extensive experimental results on the KITTI dataset prove that our framework: i) outperforms existing models for self-supervised depth estimation and ii) it is practical and couples with complex and compact models, allowing for accurate, real-time monocular depth estimation on high-end GPUs as well as on embedded systems.

# References

[1] F. Aleotti, F. Tosi, M. Poggi, and S. Mattoccia. Generative adversarial networks for unsupervised monocular depth prediction. In *15th European Conference on Computer Vision (ECCV) Workshops*, 2018. 2

[2] A. Atapour-Abarghouei and T. P. Breckon. Real-time monocular depth estimation using synthetic data with domain adaptation via image style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 18, page 1, 2018. 2

[3] V. Casser, S. Pirk, R. Mahjourian, and A. Angelova. Unsupervised learning of depth and ego-motion: A structured approach. In *Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)*, 2019. 2, 7

[4] N. Chodosh, C. Wang, and S. Lucey. Deep convolutional compressed sensing for lidar depth completion. In *Computer Vision - ACCV 2018 - 14th Asian Conference on Computer Vision, Perth, Australia, December 2-6, 2018, Revised Selected Papers, Part I*, pages 499–513, 2018. 3

[5] 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. *CoRR*, abs/1604.01685, 2016. 6, 7

[6] A. CS Kumar, S. M. Bhandarkar, and P. Mukta. Monocular depth prediction using generative adversarial networks. In *1st International Workshop on Deep Learning for Visual SLAM, (CVPR)*, 2018. 2

[7] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. *CoRR*, abs/1406.2283, 2014. 2, 5, 6, 7

[8] A. Eldesokey, M. Felsberg, and F. S. Khan. Propagating confidences through cnns for sparse data regression. In *British Machine Vision Conference 2018, BMVC 2018, Northumbria University, Newcastle, UK, September 3-6, 2018*, page 14, 2018. 3

[9] J. Engel, V. Koltun, and D. Cremers. Direct sparse odometry. Mar. 2018. 2

[10] J. Engel, T. Schöps, and D. Cremers. LSD-SLAM: Large-scale direct monocular SLAM. September 2014. 2

[11] J. Engel, J. Stückler, and D. Cremers. Large-scale direct slam with stereo cameras. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1935–1942. IEEE, 2015. 2

[12] N. Fanani, A. Sturck, M. Barnada, and R. Mester. Multimodal scale estimation for monocular visual odometry. In *IEEE Intelligent Vehicles Symposium, IV 2017, Los Angeles, CA, USA, June 11-14, 2017*, pages 1714–1721, 2017. 2, 3, 6

[13] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao. Deep ordinal regression network for monocular depth estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 2

[14] R. Garg, V. K. B. G, and I. D. Reid. Unsupervised CNN for single view depth estimation: Geometry to the rescue. *CoRR*, abs/1603.04992, 2016. 2, 4

[15] A. Geiger. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, CVPR '12, pages 3354–3361, Washington, DC, USA, 2012. IEEE Computer Society. 1, 2, 5, 6, 7

[16] C. Godard, O. Mac Aodha, and G. J. Brostow. Unsupervised monocular depth estimation with left-right consistency. *CoRR*, abs/1609.03677, 2016. 1, 2, 3, 4, 5, 6, 7, 8

[17] X. Guo, H. Li, S. Yi, J. Ren, and X. Wang. Learning monocular depth by distilling cross-domain stereo networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 484–500, 2018. 2

[18] R. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, New York, NY, USA, 2 edition, 2003. 1

[19] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 4

[20] H. Hirschmuller. Stereo processing by semiglobal matching and mutual information. *IEEE Transactions on pattern analysis and machine intelligence*, 30(2):328–341, 2008. 2, 7

[21] I. P. Howard. Perceiving in depth, vol. 1: Basic mechanisms. 2012. 1

[22] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. 6

[23] M. Klodt and A. Vedaldi. Supervising the new with the old: learning sfm from sfm. In *The European Conference on Computer Vision (ECCV)*, September 2018. 2

[24] Y. Kuznietsov, J. Stückler, and B. Leibe. Semi-supervised deep learning for monocular depth map prediction. *CoRR*, abs/1702.02706, 2017. 2

[25] L. Ladicky, J. Shi, and M. Pollefeys. Pulling things out of perspective. In *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, pages 89–96, 2014. 2

[26] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 239–248. IEEE, 2016. 2

[27] F. Liu, C. Shen, G. Lin, and I. Reid. Learning depth from single monocular images using deep convolutional neural fields. *IEEE transactions on pattern analysis and machine intelligence*, 38(10):2024–2039, 2016. 2

[28] Y. Luo, J. Ren, M. Lin, J. Pang, W. Sun, H. Li, and L. Lin. Single view stereo matching. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[29] F. Ma and S. Karaman. Sparse-to-dense: Depth prediction from sparse depth samples and a single image. In *2018 IEEE International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018*, pages 1–8, 2018. 2, 3

[30] R. Mahjourian, M. Wicke, and A. Angelova. Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2, 7

[31] R. Mur-Artal and J. D. Tardós. ORB-SLAM2: an open-source SLAM system for monocular, stereo and RGB-D cameras. *CoRR*, abs/1610.06475, 2016. 2, 3, 6

[32] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison. Dtam: Dense tracking and mapping in real-time. In *2011 international conference on computer vision*, pages 2320–2327. IEEE, 2011. 2

[33] V. Peluso, A. Cipolletta, A. Calimera, M. Poggi, F. Tosi, and S. Mattoccia. Enabling energy-efficient unsupervised monocular depth estimation on armv7-based platforms. In *Design Automation and Test in Europe (DATE)*, 2019. 2

[34] A. Pilzer, S. Lathuiliere, N. Sebe, and E. Ricci. Refine and distill: Exploiting cycle-inconsistency and knowledge distillation for unsupervised monocular depth estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 6

[35] M. Poggi, F. Aleotti, F. Tosi, and S. Mattoccia. Towards real-time unsupervised monocular depth estimation on cpu. In *IEEE/JRS Conference on Intelligent Robots and Systems (IROS)*, 2018. 2, 4, 6, 7, 8

[36] M. Poggi, F. Tosi, and S. Mattoccia. Learning monocular depth estimation with unsupervised trinocular assumptions. In *6th International Conference on 3D Vision (3DV)*, 2018. 2, 7

[37] A. Saxena, M. Sun, and A. Y. Ng. Make3d: Learning 3d scene structure from a single still image. *IEEE Trans. Pattern Anal. Mach. Intell.*, 31(5):824–840, May 2009. 2

[38] J. Shen and S.-C. S. Cheung. Layer depth denoising and completion for structured-light rgb-d cameras. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2013. 2

[39] K. Tateno, F. Tombari, I. Laina, and N. Navab. Cnn-slam: Real-time dense monocular slam with learned depth prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6243–6252, 2017. 2

[40] F. Tosi, F. Aleotti, M. Poggi, and S. Mattoccia. Learning monocular depth estimation infusing traditional stereo knowledge. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2, 4, 7

[41] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox, and A. Geiger. Sparsity invariant cnns. *CoRR*, abs/1708.06500, 2017. 2, 3

[42] C. Wang, J. Miguel Buenaposada, R. Zhu, and S. Lucey. Learning depth from monocular videos using direct methods. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2, 7

[43] X. Wang, H. Zhang, X. Yin, M. Du, and Q. Chen. Monocular visual odometry scale recovery using geometrical constraint. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 988–995, May 2018. 2, 3, 6

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

[45] D. Xu, W. Wang, H. Tang, H. Liu, N. Sebe, and E. Ricci. Structured attention guided convolutional neural fields for monocular depth estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[46] N. Yang, R. Wang, J. Stückler, and D. Cremers. Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry. In *European Conference on Computer Vision*, pages 835–852. Springer, 2018. 2, 5, 6

[47] X. Yin, X. Wang, X. Du, and Q. Chen. Scale recovery for monocular visual odometry using depth estimated with deep convolutional neural fields. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5870–5878, 2017. 2

[48] Z. Yin and J. Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2, 7

[49] P. Zama Ramirez, M. Poggi, F. Tosi, S. Mattoccia, and L. Di Stefano. Geometry meets semantic for semi-supervised monocular depth estimation. In *14th Asian Conference on Computer Vision (ACCV)*, 2018. 2

[50] H. Zhan, R. Garg, C. Saroj Weerasekera, K. Li, H. Agarwal, and I. Reid. Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2

[51] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe. Unsupervised learning of depth and ego-motion from video. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 1, 2, 7

[52] Y. Zou, Z. Luo, and J.-B. Huang. Df-net: Unsupervised joint learning of depth and flow using cross-task consistency. In *European Conference on Computer Vision*, 2018. 2, 7