




Pose Constraints for Consistent Self-supervised Monocular Depth and Ego-Motion

Zeeshan Khan Suri^(✉) 

DENSO ADAS Engineering Services GmbH, Kemptener Str. 99,
88131 Lindau, Germany
z.suri@eu.denso.com

Abstract. Self-supervised monocular depth estimation approaches suffer not only from scale ambiguity but also infer temporally inconsistent depth maps w.r.t. scale. While disambiguating scale during training is not possible without some kind of ground truth supervision, having scale consistent depth predictions would make it possible to calculate scale once during inference as a post-processing step and use it over-time. With this as a goal, a set of temporal consistency losses that minimize pose inconsistencies over time are introduced. Evaluations show that introducing these constraints not only reduces depth inconsistencies but also improves the baseline performance of depth and ego-motion prediction.

Keywords: self-supervision · 3d-reconstruction · depth-estimation

1 Introduction

Depth is so essential for perceiving, understanding and navigating the 3D world around us that biological animals have evolved to possess redundant apparatus for perceiving it. Animals are able to infer relative depth of the perceived scene even with a single eye [41]. Although the loss of dimensionality from projecting a 3D scene on a 2D plane cannot be completely recovered, an image consists of many useful monocular cues, such as relative sizes, texture gradient, linear perspective, contrast differences. Image sequences contain additional motion cues such as occlusion, motion parallax [3]. These cues impose constraints on the possible combinations of depth of the underlying scene. Such cues can be implicitly learnt by computational models in a supervised [8, 36], and in a self-supervised [10, 12, 59] manner, making it possible to infer relative depth from a single image. During training, the self-supervised methods rely on motion cues coming from image sequences. They simultaneously estimate relative pose between two consecutive frames of a sequence and their individual depth such that by warping one frame onto the other using the depth, relative pose and camera intrinsics, the other frame can be synthesized. The underlying assumption is of photometric constancy of neighboring frames of a sequence, and thus a loss between the warped neighboring image and the true one can back-propagate through the depth and pose networks.

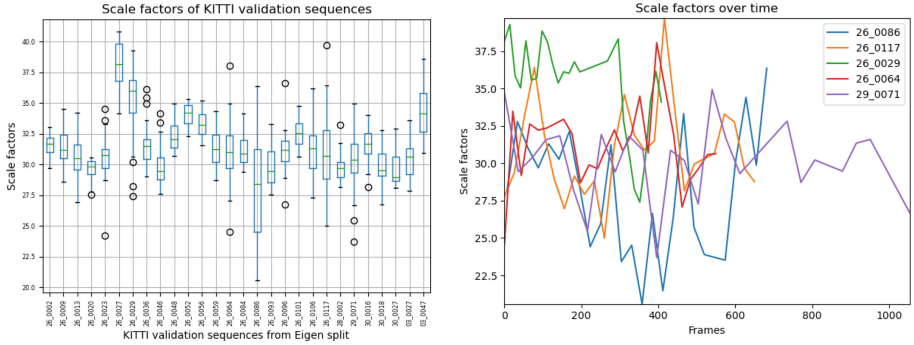


Fig. 1. Variation of scale factors $\frac{\text{median}(D_{\text{true}})}{\text{median}(D_{\text{pred}})}$ per frame, within each KITTI validation sequence is shown (**Left**): as a box-plot, and (**Right**): over time.

The 2D image pixels might have been projected from infinite number of 3D points, thus making the inverse problem of recovering the depth dimension from the 2D image ill-posed. While cues constrain the solution depth map to have an underlying structure, the estimated depth can only be relative in nature and not absolute, since any scalar multiple of the estimated depth map could also be equally optimal. Furthermore, these self-supervised models do not necessarily learn our standard units of measuring distance, but rather predict depth in their own. Even for humans these units, such as metric, imperial, US customary, tend to vary in usage, for, after all, these units and the scalar factors relating them are a human construct. Such scale ambiguity is not usually a problem, since one can calculate the scale factor relating the model’s depth to that of metric or imperial as a post-processing step and convert one to the other. Unfortunately, the depth estimates from self-supervised models are not just scale ambiguous but also temporally scale inconsistent, *i.e.* the scale of one frame’s depth map is different to that of the neighboring frame’s. This variation of the scale factors is shown in Fig. 1 as a box-plot for each KITTI validation sequence and also over-time for the top 5 varying sequences. Here, the scale factor of a frame is the ratio of the metric LIDAR depth to that of the predicted depth $= \frac{\text{median}(D_{\text{true}})}{\text{median}(D_{\text{pred}})}$.

Due to this variation, it is not accurate to analytically determine the scale factor once and use it later, and makes it necessary to calculate the scale in each frame, making it infeasible for applications without availability of some kind of ground-truth, as in case of online videos. In fact, most current methods artificially scale their outputs depth maps per-frame, using the LIDAR ground truth, as a post-processing step before calculating the error at test time.

This work hypothesises that the scale consistency problem is due to the lack of proper temporal constraints, *i.e.* the scale is independently ambiguous for each frame, and the training pipeline finds a scale that is locally optimal for that frame, regardless if it agrees with the scale of the neighboring frames. While some of the recent works on this issue impose consistency in depth, this work instead

explores additional ways of imposing temporal constraints in an unsupervised manner, *i.e.* without the need of additional ground truth supervision.

2 Related Work

This section explores different ways the scale ambiguity and the scale consistency problem has been tackled in the literature. Also losses similar to ours in related tasks are referred. Temporal data such as video has temporal correlations and has a low probability to abruptly change within short intervals. Temporal consistency exploits this information flow and has been utilized for various video applications [2, 22, 28], including supervised depth estimation [55], style transfer [5], video completion [15] segmentation [56] to capture temporal correlations. Temporal consistency is also the key component of the recent self-supervised monocular depth estimation approaches [19, 59], where the photometric constancy among consecutive frames is assumed and differences after reprojection are minimized. But, since the depth network in these methods is monocular and sees only one frame at a time, the predicted depth maps are not consistent over time.

2.1 Scale-Disambiguation/Consistency-Enforcement via Supervision

Methods with some kind of ground truth supervision are able to enforce scale consistency and do not have the same scale inconsistency issues, for example, the scale factor in stereo methods comes from the relation of disparity and depth as $\text{depth} = \frac{\text{focal length} \times \text{baseline}}{\text{disparity}}$, and is constant. Following this, Roussel et al. [34] first pre-train on stereo data and then fine-tune on monocular data, to show that doing so retains the scale learnt from the stereo pretraining. GPS or ground-truth pose data has also been used to disambiguate scale [4, 12] via enforcing the pose network’s output to match the pose coming from sensors.

Hand-picked features based depth estimation methods such as SLAM [31] and Structure from Motion (SfM) [37] have been used [26, 40] for a source of supervision, to transfer wide baseline symmetric depth and sparse long-term depth consistencies from it to the depth estimation network. In a similar fashion, ideas from the Visual Odometry (VO), such as epipolar geometry and bundle adjustment, are incorporated [54, 58], to independently compute correspondences and triangulate them to produce sparse depth, which acts as additional supervision to the network’s depth prediction.

Via the Plane Assumption. Following Kitt *et al.* [20]’s scale recovery for VO, many recent works [29, 52] make use of the following assumptions: a) most automotive cameras are rigidly mounted in a fixed position and at constant orientation with respect to the road, b) the roll and pitch movement of the vehicle have negligible effect on its position and orientation, c) most urban streets may be assumed to be approximately planar in the vicinity of the vehicle. These assumptions allow them to use camera extrinsic parameters, in particular the camera height and compare it with the estimated height by fitting plane on the road, to recover the scale as a post-processing step. [44, 50] also use camera height

for scale but incorporate it within training. These methods rely on heuristics that a flat road plane is visible in the area of interest and that the camera position and orientation remain constant over time, which are often not realistic.

Methods that use some kind of supervision for disambiguating depth at each frame, also inherently enforce temporal consistency. Since the supervision is temporally consistent, the predicted depth maps are also made temporally consistent as an unintended consequence.

2.2 Self-supervised Temporal Consistency

Recurrent neural networks [32, 46] have been used to implicitly model multiple frames inputs. Having multiple frames as input directly to the depth network causes holes at regions with moving objects [49], and need to be corrected by an additional single input network. 3D geometric constraints [6, 27, 45] were proposed that penalize the euclidean distances between the reconstructed point clouds of two consecutive frames, after transforming one to the other. Similar geometric constraints on the depth maps were proposed [1, 25, 35] which minimize the inconsistency of the estimated disparity maps of two consecutive frames, after warping one onto the other. Our work lies in this category and differs in the fact that we propose complementary constraints on the pose, which can, in principle, be added and used together with the other temporal constraints.

2.3 Similar Constraints in Literature

Constraints similar to the ones we impose via a loss are found across various computer vision applications.

Forward-Backward Consistency. Li et al. [24] propose a forward-backward pose consistency loss but in a stereo setting, where, the pose from one frame of the left camera to its neighboring frame should be identical to the pose between the same frame of the right camera to its neighboring frame. Based on Narayanan *et al.* [39]’s forward-backward optical flow consistency assumption, [30] propose a loss in their optical flow prediction network. This was adopted in the context of monocular depth and ego-motion [46, 53], where the flow caused by rigid ego-motion estimation is computed and the forward-backward inconsistency of the rigid flow is minimized. Sheng *et al.* [38] propose a loss that minimizes the forward-backward inconsistencies in the bi-directional warping fields generated from the rigid ego-motion. [54, 62] train an optical flow network, in addition to the depth and ego-motion networks for per-pixel dense 2D pose between two consecutive frames and adopt the same forward-backward consistency loss for their optical flow. Li et al. [23] propose a forward-backward loss directly on pose estimates. Our forward-backward loss is similar to theirs.

Cycle Consistency. The forward-backward consistency is an instance of the idea of cycle-consistency which has been applied to a wide range of computer vision tasks [14, 16, 47, 60]. Related to our problem of interest, pose cycle consistency was used in the context of Visual Odometry [18, 61]. Ruhkamp *et al.* [35]

propose a strategy to detect and mask inconsistent regions such as occlusions, in neighboring depth maps via a cycle consistency of the reprojected RGB frames.

3 Method

Given video sequences captured from a camera with known intrinsic parameters K , the objective is to learn a depth network model $f_D : \mathbb{R}^{3 \times H \times W} \rightarrow \mathbb{R}^{H \times W}$ that maps an RGB frame at time t , $I_t \in \mathbb{R}^{3 \times H \times W}$ to the depth of its underlying scene $D_t \in \mathbb{R}^{H \times W}$, and an ego-motion network model $f_E : \mathbb{R}^{2(3 \times H \times W)} \rightarrow \mathfrak{se}(3)$ that take in two consecutive RGB frames $\{I_t, I_{t+n}$, typically, $n \in \{-1, 1\}\}$ and output the 6°C-of-Freedom (DOF) rigid transformation $\mathbf{T}_t^{t+n} = (\mathbf{e}^T, \mathbf{t}^T)^T \in \mathfrak{se}(3)$, and, $\begin{pmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{pmatrix} \in \text{SE}(3)$ between them, where H, W are the image height and width and, $\mathbf{e} \in \mathfrak{so}(3)$, $\mathbf{t} \in \mathbb{R}^3$, $\mathbf{R} \in \text{SO}(3)$ represent the axis-angle, translation and the rotation matrix respectively.

The depth and the ego-motion are jointly trained with the key assumption of photometric constancy, *i.e.* a view reprojected with the camera intrinsics, depth and the camera motion onto its neighboring view, reconstructs the neighboring view. The reprojection of a view onto the neighbor is as follows

$$[\hat{u}, \hat{v}, 1]^T \sim K \mathbf{R}_t^{t+n} D_t^{ij} K^{-1} [u, v, 1]^T + K \mathbf{t}_t^{t+n}, \quad (1)$$

where $[u, v, 1]$ and $[\hat{u}, \hat{v}, 1]$ denote the homogeneous pixel coordinates of view at time t and that of its neighboring view, D_t^{ij} is the pixel's corresponding depth and $\mathbf{R}_t^{t+n}, \mathbf{t}_t^{t+n}$ denote the rotation and translation from the view at time t to that of the neighboring view at time $t+n$ respectively. Using the reprojection Eq. (1), the view at time t is synthesized from the neighboring view as

$$\hat{I}_{t+n \rightarrow t}[u, v] = I_{t+n} \langle [\hat{u}, \hat{v}] \rangle, \quad (2)$$

where $\langle \rangle$ denotes the sampling operator. A robust photometric loss [11, 57] between the synthesized RGB view $\hat{I}_{t+n \rightarrow t}$ to the original one I_t , as in Eq. (2), is minimized, thereby updating and correcting the depth and ego-motions model weights via backpropagation.

$$\mathcal{L}_p(I_t, \hat{I}_{t+n \rightarrow t}) = \frac{\alpha}{2} (1 - \text{SSIM}(I_t, \hat{I}_{t+n \rightarrow t})) + (1 - \alpha) \|I_t - \hat{I}_{t+n \rightarrow t}\|_1, \quad (3)$$

where $\alpha = 0.85$ and SSIM denotes the structural similarity loss [48]. Following Godard *et al.* [11] an additional edge δ_x, δ_y aware smoothness term regularizes the predicted depth:

$$\mathcal{L}_s = \sum_{uv} |\delta_x \bar{D}_t^{uv}| e^{-|\delta_x I_t^{uv}|} + |\delta_y \bar{D}_t^{uv}| e^{-|\delta_y I_t^{uv}|}, \quad (4)$$

where \bar{D}_t^{uv} is the mean normalized depth at pixel $[u, v]$. We establish a strong baseline by following the best practices from Monodepth2 [10]. We also use the minimum reprojection error $\min_n \mathcal{L}_p(I_t, \hat{I}_{t+n \rightarrow t})$, over pairs of both neighboring

frames to deal with occlusions and auto-masking to disregard temporally static pixels. These losses are minimized over all pixels in the training set over four resolution scales. The readers are referred to Monodepth2 [10] for more details.

3.1 Pose Constraints

The depth and the pose networks are tightly coupled. We propose to establish temporal consistency in the predicted depth through the pose network’s estimates of the ego-motion between two input frames. Specifically, we propose three constraints that do not add any additional assumptions, but are expected to be met explicitly. Let $T = \begin{pmatrix} \mathbf{R}_T & \mathbf{t}_T \\ \mathbf{0} & 1 \end{pmatrix} \in \text{SE}(3)$ and $P = \begin{pmatrix} \mathbf{R}_P & \mathbf{t}_P \\ \mathbf{0} & 1 \end{pmatrix} \in \text{SE}(3)$ be two 6-DOF rigid poses that represent the same transformation, we define the distance metric $d(T, P) : \text{SE}(3) \times \text{SE}(3) \rightarrow \mathbb{R}^+$, between them as

$$d(T, P) = d(\mathbf{R}_T, \mathbf{R}_P) + d(\mathbf{t}_T, \mathbf{t}_P) = \left\| 1 - \left(\frac{\text{tr}(\mathbf{R}_T \mathbf{R}_P^T) - 1}{2} \right) \right\|_1 + \|\mathbf{t}_T - \mathbf{t}_P\|_1 \quad (5)$$

where $\theta = \cos^{-1} \left(\frac{\text{tr}(\mathbf{R}_T \mathbf{R}_P^T) - 1}{2} \right) \in [0, \pi]$ is the angle of the relative rotation $\mathbf{R}_T \mathbf{R}_P^T$, and $1 - \cos(\theta)$ corresponds to the geodesic distance on the 3D manifold of rotation matrices [17, 43]. Throughout the proposed constraints, we use Eq. (5) as the distance metric between $\text{SE}(3)$ transformations.

Forward-Backward Pose Consistency. The ego-motion between frames at time t and $t + n$ should be the inverse of the ego-motion between frames at time $t + n$ and t , *i.e.* $\mathbf{T}_{t+n}^t \stackrel{!}{=} \mathbf{T}_t^{t+n-1}$. This is characterized by minimizing the following loss

$$\mathcal{L}_{\mathbf{T}_{fb}} = d(\mathbf{T}_{t+n}^t, \mathbf{T}_t^{t+n-1}), \quad (6)$$

where $\mathbf{T}^{-1} = \begin{pmatrix} \mathbf{R}^T & -\mathbf{t} \\ \mathbf{0} & 1 \end{pmatrix}$

Identity Pose Consistency. Self-supervised monocular depth estimation operates on the underlying assumption of moving camera and previous works address this by proposing masking strategies to filter out pixels with no overall motion, either due to static camera or moving objects [10]. This work, in addition, forces the pose network to explicitly inspect static images and estimate no relative ego-motion. This is done by giving the pose network same frames and having a loss that minimizes any ego-motion $\mathbf{T}_t^t = \begin{pmatrix} \mathbf{e}_t^t \\ \mathbf{t}_t^t \end{pmatrix} \in \mathfrak{se}(3)$. The loss is as follows

$$\mathcal{L}_{\mathbf{T}_{id}} = \|\mathbf{T}_t^t\|_1 = \|\mathbf{e}_t^t\|_1 + \|\mathbf{t}_t^t\|_1 \quad (7)$$

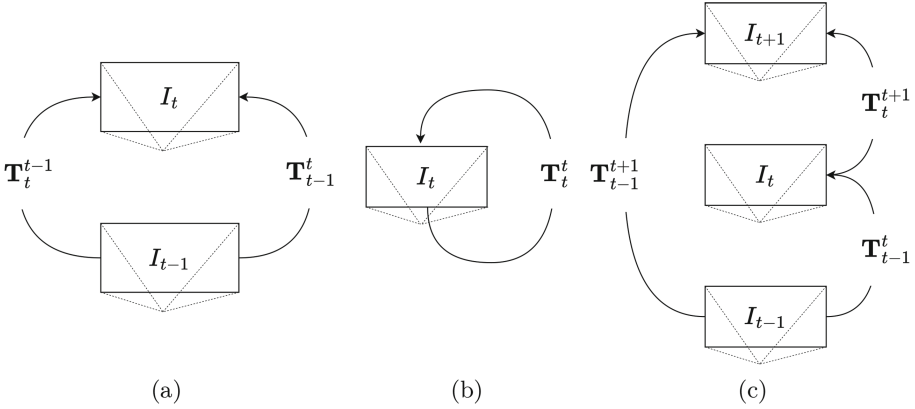


Fig. 2. Illustration of the proposed: (a) forward-backward, (b) identity and (c) cycle, pose constraints

Cycle Pose Consistency. Cyclic pose consistency states that the poses need to be consistent in a cyclic manner, *i.e.* the combined pose from $t-n$ to $t+n$ via t , should be the same as the direct ego-motion between them, $\mathbf{T}_{t-n}^{t+n} \stackrel{!}{=} \mathbf{T}_t^{t+n} \times \mathbf{T}_{t-n}^t$, as illustrated in Fig. 2(c). The loss is parameterized as follows

$$\mathcal{L}_{\mathbf{T}_{cyc}} = d(\mathbf{T}_{t-n}^{t+n}, \mathbf{T}_t^{t+n} \times \mathbf{T}_{t-n}^t) \quad (8)$$

4 Experiments

To maintain consistency and comparability, the exact experimental setup of baseline [10] is maintained. Following the established protocols by Eigen et al. [7], we train and test on Zhou *et al.* [59]’s splits of the KITTI [9] raw dataset, with a widely-used pre-processing to remove static frames [59] and cap the depth at 80m. We evaluate the depth model using popularly used metrics from Eigen et al. [7]. We also consider the improved ground truth depth maps [42] for evaluation, which uses stereo information and 5 consecutive frames to accumulate LiDAR points to handle moving objects, resulting in high quality depth maps. We weight our loss with a weight of 0.1 when addition to the final objective.

4.1 Temporal Scale Consistency of Predicted Depth

The main goal of this work is to minimize scale inconsistencies. We introduce a simple metric to measure the scale consistency across consecutive frames

Consistency Metric. We hypothesize that in practice, it would be possible to calculate the actual scale once and use it for the rest of the sequence. Thus, the actual value of the scale does not matter but only its normalized variation w.r.t.

Table 1. Depth evaluation with each of our constraints applied individually. The second column shows the coefficient of correlation (normalized standard deviation) of scales as a measure of consistency. The top section uses ground-truth scaling in order to make the depth metrics comparable. In the bottom section, a per-sequence (in contrast to per-frame) median scaling is used. Our constraints show reduction in inconsistencies, while also slightly improving the depth.

Constraint	Lower is better					Accuracy (Higher is better)		
	$\frac{\sigma(\text{scales})}{\mu(\text{scales})}$	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
–	0.096	0.116	0.896	4.915	0.194	0.872	0.957	0.981
$\mathcal{L}_{T_{id}}$	0.090	0.116	0.897	4.823	0.193	0.874	0.959	0.981
$\mathcal{L}_{T_{fb}}$	0.088	0.113	0.867	4.796	0.191	0.878	0.960	0.981
$\mathcal{L}_{T_{cyc}}$	0.090	0.113	0.851	4.809	0.191	0.877	0.959	0.981
–	–	0.120	0.924	4.962	0.198	0.860	0.956	0.980
$\mathcal{L}_{T_{fb}}$	-	0.116	0.885	4.855	0.195	0.867	0.958	0.981

time. We define the coefficient of correlation of the computed scale factors as $\frac{\sigma(\text{scales})}{\mu(\text{scales})}$, where $\text{scale} = \frac{\text{median}(D_{\text{true}})}{\text{median}(D_{\text{pred}})}$, as a measure of temporal scale consistency in depth.

Table 1 compares our proposed constraints on the standard depth metrics as well as the proposed consistency metric. Additionally, for each sequence in the KITTI test set, we calculate the scale as the median of all per-frame scales and use it throughout the sequence, resulting in the evaluation shown in Table 1’s bottom two rows. Our constraints show reduction in inconsistencies, while also slightly improving the depth.

4.2 Depth Evaluation

Table 2 compares our best model (with $\mathcal{L}_{T_{cyc}}$) with works which tackle the scale ambiguity/scale consistency problem. Works which make use of some kind of ground truth supervision, shown in the top section of Table 2, disambiguate scale as a byproduct and, as a result, are automatically scale consistent. Different types of supervision signals have been used as mentioned in the *Supervision* column. In the bottom section, works tackling the problem without additional supervision are mentioned. Our work belongs in this category.

Table 2 shows that our proposed pose consistency constraint gives improved performance with respect to baseline [10], although this was not our goal, showing that temporal consistency is also important for the accuracy

4.3 Ego-motion Evaluation

Following Zhou *et al.* [59]’s protocols, we also evaluate our ego-motion network. We train on KITTI odometry splits sequences 0–8 and test on sequences 9 and 10. Similar to the related methods, we compute the absolute trajectory error (ATE) averaged over all overlapping 5-frame snippets. Since our pose network

Table 2. Comparison of our results with works that focus on temporal consistency, on the test set of KITTI [9]’s Eigen [7] split. Works in the bottom section use monocular (M) supervision, while those in the top section use additional supervision: D for LIDAR depths, h for camera extrinsics (height), v for velocity (GT pose) and SfM/SLAM for depth hints from classical methods. For fair comparison, only the contributions tackling the inconsistency problem are compared against. The best and second best methods are in bold and underlined respectively.

Method	Supervision	Error (Lower is better)				Accuracy (Higher is better)		
		Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Sparse-to-Cont [33]	M+D	0.118	0.630	4.520	0.209	0.898	0.966	0.985
DNet [52]	M+h	0.113	0.864	4.812	0.191	0.877	0.960	0.981
pRGBD-Refined [40]	M+SLAM	0.113	0.793	4.655	0.188	0.874	0.960	0.983
Kuznetsov <i>et al.</i> [21]	M+D	0.113	0.741	4.621	0.189	0.862	0.960	0.986
TrainFlow [58]	M+SfM	0.113	0.704	4.581	0.184	0.871	0.961	0.984
G2S R50 [4]	M+v	0.112	0.894	4.852	0.192	0.877	0.958	0.981
PackNet-SfM [12] (ResNet18)	M+v	0.111	0.829	4.788	0.199	0.864	0.954	0.980
VA-Depth [51]	M+h	0.112	0.864	4.804	0.190	0.878	0.959	0.982
vid2depth [27]	M	0.163	1.240	6.220	0.250	0.762	0.916	0.968
DF-Net [62]	M	0.150	1.124	5.507	0.223	0.806	0.933	0.973
Sheng <i>et al.</i> [38]	M	0.139	1.021	5.418	0.209	0.803	0.937	0.976
Li <i>et al.</i> [23]	M	0.130	0.950	5.138	0.209	0.843	0.948	0.978
SC-SfMLearner [1]	M	0.119	0.857	4.950	0.197	0.863	0.957	<u>0.981</u>
TC-Depth [35] (only \mathcal{L}_{geo})	M	<u>0.113</u>	0.904	<u>4.773</u>	0.193	<u>0.877</u>	<u>0.959</u>	0.980
Wang <i>et al.</i> [45]	M	0.109	0.779	4.641	0.186	0.883	0.962	0.982
Baseline (Monodepth2) [10]	M	0.115	0.903	4.863	0.193	<u>0.877</u>	<u>0.959</u>	<u>0.981</u>
Ours	M	<u>0.113</u>	<u>0.851</u>	4.809	<u>0.191</u>	<u>0.877</u>	<u>0.959</u>	<u>0.981</u>

only takes 2 frames as input, we aggregate the relative poses to create 5-frame trajectories. Table 3 summarizes our improvements with respect to our baseline [10]. Our proposed loss constrains the solution space, thereby making a better optimum easier to achieve.

Table 3. Ego-motion estimation results: average absolute trajectory error, and standard deviation, in meters, on KITTI Odometry dataset [9]. Trained on Seq. 0-8 and tested on Seq. 9 and 10. Our pose constraints show improvement with respect to baseline [10].

	Seq. 09	Seq. 10
DF-Net [62]	0.017 \pm 0.007	0.015 \pm 0.009
Wang <i>et al.</i> [45]	0.014 \pm 0.008	0.014 \pm 0.010
Baseline [10]	0.017 \pm 0.008	0.015 \pm 0.010
Ours ($\mathcal{L}_{\mathcal{T}_{cyc}}$)	0.016 \pm 0.008	0.014 \pm 0.010

5 Conclusion and Future Work

Self-supervised monocular depth and ego-motion estimation methods suffer not only from scale ambiguity but also from scale inconsistency. While including

some kind of ground truth (GT) supervision not only disambiguates scale but also enforces scale consistency, it is not always plausible to have access to accurate ground-truth information. We propose ego-motion constraints that do not require any additional GT supervision. Via experimentation, we show that our proposed constraints not only decrease the inconsistencies but also improve the depth's and ego-motion's performance compared to baseline.

Our constraints do not aim to compete but complement the ones used in literature, for example, SC-SfMLearner [1]'s depth consistency. We have looked at how individual constraint performs. We leave the effect of their interactions with one another and constraints used by previous works for the future. We also do not loosen the static scene assumption. If the image motion is dominated by moving objects, it is indeed difficult to estimate the true ego-motion caused just by cameras. It would be interesting to generalize these constraints with denser motion models, such as Li et al. [23]'s piece-wise rigid flows considering dynamic objects or dense per-pixel optical flows [13].

References

1. Bian, J.W., et al.: Unsupervised scale-consistent depth learning from video. *Int. J. Comput. Vision (IJCV)* (2021)
2. Bonneel, N., Tompkin, J., Sunkavalli, K., Sun, D., Paris, S., Pfister, H.: Blind video temporal consistency. *ACM Trans. Graph.* 34(6) (2015). <https://doi.org/10.1145/2816795.281810>
3. Burton, H.E.: The optics of euclid1. *J. Opt. Soc. Am.* **35**(5), 357–372 (1945)
4. Chawla, H., Varma, A., Arani, E., Zonooz, B.: Multimodal scale consistency and awareness for monocular self-supervised depth estimation. In: 2021 IEEE International Conference on Robotics and Automation (ICRA). IEEE (in press) (2021)
5. Chen, D., Liao, J., Yuan, L., Yu, N., Hua, G.: Coherent online video style transfer. In: 2017 IEEE International Conference on Computer Vision (ICCV), pp. 1114–1123 (2017). <https://doi.org/10.1109/ICCV.2017.126>
6. Chen, Y., Schmid, C., Sminchisescu, C.: Self-supervised learning with geometric constraints in monocular video: Connecting flow, depth, and camera. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 7062–7071 (2019). <https://doi.org/10.1109/ICCV.2019.00716>
7. Eigen, D., Puhrsch, C., Fergus, R.: Depth map prediction from a single image using a multi-scale deep network. In: Proceedings of the 27th International Conference on Neural Information Processing Systems, NIPS 2014, vol. 2. pp. 2366–2374. MIT Press, Cambridge (2014)
8. Fu, H., Gong, M., Wang, C., Batmanghelich, K., Tao, D.: Deep ordinal regression network for monocular depth estimation. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2002–2011 (2018). <https://doi.org/10.1109/CVPR.2018.00214>
9. 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 (2012). <https://doi.org/10.1109/CVPR.2012.6248074>
10. Godard, C., Mac Aodha, O., Firman, M., Brostow, G.J.: Digging into self-supervised monocular depth prediction (October 2019)

11. Godard, C., Aodha, O.M., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6602–6611 (2017). <https://doi.org/10.1109/CVPR.2017.699>
12. Guizilini, V., Ambrus, R., Pillai, S., Raventos, A., Gaidon, A.: 3d packing for self-supervised monocular depth estimation. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
13. Guizilini, V., Lee, K.H., Ambrus, R., Gaidon, A.: Learning optical flow, depth, and scene flow without real-world labels. *IEEE Robotics Automation Lett.* **7**(2), 3491–3498 (2022). <https://doi.org/10.1109/LRA.2022.3145057>
14. Hoffman, J., et al.: CyCADA: Cycle-consistent adversarial domain adaptation. In: Dy, J., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 80, pp. 1989–1998. PMLR (10–15 Jul 2018), <https://proceedings.mlr.press/v80/hoffman18a.html>
15. Huang, J.B., Kang, S.B., Ahuja, N., Kopf, J.: Temporally coherent completion of dynamic video. *ACM Trans. Graph.* **35**(6) (2016). <https://doi.org/10.1145/2980179.2982398>
16. Huang, Q.X., Guibas, L.: Consistent shape maps via semidefinite programming. *Computer Graphics Forum* **32**(5), 177–186 (2013)
17. Huynh, D.Q.: Metrics for 3d rotations: Comparison and analysis. *J. Mathemat. Imaging Vision* **35**, 155–164 (2009)
18. Iyer, G., Krishna Murthy, J., Gupta, G., Krishna, M., Paull, L.: Geometric consistency for self-supervised end-to-end visual odometry. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp. 267–275 (2018)
19. Karsch, K., Liu, C., Kang, S.B.: Depth transfer: Depth extraction from video using non-parametric sampling. *IEEE Trans. Pattern Anal. Mach. Intell.* **36**(11), 2144–2158 (2014). <https://doi.org/10.1109/TPAMI.2014.2316835>
20. Kitt, B., Rehder, J., Chambers, A., Schönbein, M., Lategahn, H., Singh, S.: Monocular visual odometry using a planar road model to solve scale ambiguity (2011)
21. Kuznetsov, Y., Stückler, J., Leibe, B.: Semi-supervised deep learning for monocular depth map prediction. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2215–2223 (2017). <https://doi.org/10.1109/CVPR.2017.238>
22. Lang, M., Wang, O., Aydin, T., Smolic, A., Gross, M.: Practical temporal consistency for image-based graphics applications. *ACM Trans. Graph.* **31**(4), 34:1–34:8 (2012). <https://doi.org/10.1145/2185520.2185530>
23. Li, H., Gordon, A., Zhao, H., Casser, V., Angelova, A.: Unsupervised monocular depth learning in dynamic scenes. In: Conference on Robot Learning, pp. 1908–1917. PMLR (2021)
24. 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 (2018). <https://doi.org/10.1109/ICRA.2018.8461251>
25. Li, S., Luo, Y., Zhu, Y., Zhao, X., Li, Y., Shan, Y.: Enforcing temporal consistency in video depth estimation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 1145–1154 (October 2021)
26. Luo, X., Huang, J.B., Szeliski, R., Matzen, K., Kopf, J.: Consistent video depth estimation. *ACM Trans. Graph.* **39**(4) (2020). <https://doi.org/10.1145/3386569.3392377>

27. Mahjourian, R., Wicke, M., Angelova, A.: Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In: CVPR (2018)
28. Mallya, A., Wang, T.-C., Saprà, K., Liu, M.-Y.: World-consistent video-to-video synthesis. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.-M. (eds.) ECCV 2020. LNCS, vol. 12353, pp. 359–378. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-58598-3_22
29. McCraith, R., Neumann, L., Vedaldi, A.: Calibrating self-supervised monocular depth estimation. In: Machine Learning for Autonomous Driving Workshop at the 34th Conference on Neural Information Processing Systems (NeurIPS 2020) (2020)
30. Meister, S., Hur, J., Roth, S.: Unflow: Unsupervised learning of optical flow with a bidirectional census loss. In: Proceedings of the AAAI Conference On Artificial Intelligence, vol. 32 (2018)
31. Mur-Artal, R., Tardós, J.D.: Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE Trans. Rob.* **33**(5), 1255–1262 (2017). <https://doi.org/10.1109/TRO.2017.2705103>
32. Patil, V., Van Gansbeke, W., Dai, D., Van Gool, L.: Don't forget the past: Recurrent depth estimation from monocular video. *IEEE Robotics Automation Lett.* **5**(4), 6813–6820 (2020). <https://doi.org/10.1109/LRA.2020.3017478>
33. Rosa, N.d.S., Guizilini, V., Grassi, V.: Sparse-to-continuous: Enhancing monocular depth estimation using occupancy maps. In: 2019 19th International Conference on Advanced Robotics (ICAR), pp. 793–800 (2019). <https://doi.org/10.1109/ICAR46387.2019.8981652>
34. Roussel, T., Eycken, L.V., Tuytelaars, T.: Monocular depth estimation in new environments with absolute scale. In: 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 1735–1741 (2019). <https://doi.org/10.1109/IROS40897.2019.8967677>
35. Ruhkamp, P., Gao, D., Chen, H., Navab, N., Busam, B.: Attention meets geometry: Geometry guided spatial-temporal attention for consistent self-supervised monocular depth estimation. In: 2021 International Conference on 3D Vision (3DV), pp. 837–847 (2021). <https://doi.org/10.1109/3DV53792.2021.00092>
36. Saxena, A., Sun, M., Ng, A.Y.: Make3d: Learning 3d scene structure from a single still image. *IEEE Trans. Pattern Anal. Mach. Intell.* **31**(5), 824–840 (2009). <https://doi.org/10.1109/TPAMI.2008.132>
37. Schönberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4104–4113 (2016). <https://doi.org/10.1109/CVPR.2016.445>
38. Sheng, L., Xu, D., Ouyang, W., Wang, X.: Unsupervised collaborative learning of keyframe detection and visual odometry towards monocular deep slam. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 4301–4310 (2019). <https://doi.org/10.1109/ICCV.2019.00440>
39. Sundaram, N., Brox, T., Keutzer, K.: Dense point trajectories by gpu-accelerated large displacement optical flow. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010. LNCS, vol. 6311, pp. 438–451. Springer, Heidelberg (2010). https://doi.org/10.1007/978-3-642-15549-9_32
40. Tiwari, L., Ji, P., Tran, Q.-H., Zhuang, B., Anand, S., Chandraker, M.: Pseudo RGB-D for self-improving monocular SLAM and depth prediction. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.-M. (eds.) ECCV 2020. LNCS, vol. 12356, pp. 437–455. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-58621-8_26
41. Trychin, S., Walk, R.D.: A study of the depth perception of monocular hooded rats on the visual cliff. *Psychonomic Sci.*, 53–54 (1964). <https://doi.org/10.3758/BF03342786>

42. Uhrig, J., Schneider, N., Schneider, L., Franke, U., Brox, T., Geiger, A.: Sparsity invariant cnns. In: 2017 International Conference on 3D Vision (3DV), pp. 11–20 (2017). <https://doi.org/10.1109/3DV.2017.00012>
43. van der, V.: (<https://math.stackexchange.com/users/76466/kwin-van-der-veen>), K.: How to compare two rotations represented by axis-angle rotation vectors? Mathematics Stack Exchange, <https://math.stackexchange.com/q/4001635>
44. Wagstaff, B., Kelly, J.: Self-supervised scale recovery for monocular depth and egomotion estimation. In: 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 2620–2627 (2021). <https://doi.org/10.1109/IROS51168.2021.9635938>
45. Wang, L., Wang, Y., Wang, L., Zhan, Y., Wang, Y., Lu, H.: Can scale-consistent monocular depth be learned in a self-supervised scale-invariant manner? In: 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 12707–12716 (2021). <https://doi.org/10.1109/ICCV48922.2021.01249>
46. Wang, R., Pizer, S.M., Frahm, J.M.: Recurrent neural network for (un-) supervised learning of monocular video visual odometry and depth. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5555–5564 (2019)
47. Wang, X., Jabri, A., Efros, A.A.: Learning correspondence from the cycle-consistency of time. In: CVPR (2019)
48. Wang, Z., Bovik, A., Sheikh, H., Simoncelli, E.: Image quality assessment: from error visibility to structural similarity. *IEEE Trans. Image Process.* **13**(4), 600–612 (2004). <https://doi.org/10.1109/TIP.2003.819861>
49. Watson, J., Aodha, O.M., Prisacariu, V., Brostow, G., Firman, M.: the temporal opportunist: self-supervised multi-frame monocular depth. In: Computer Vision and Pattern Recognition (CVPR) (2021)
50. Xiang, J., Wang, Y., An, L., Liu, H., Wang, Z., Liu, J.: Visual attention-based self-supervised absolute depth estimation using geometric priors in autonomous driving. arXiv preprint [arXiv:2205.08780](https://arxiv.org/abs/2205.08780) (2022)
51. Xiang, J., Wang, Y., An, L., Liu, H., Wang, Z., Liu, J.: Visual attention-based self-supervised absolute depth estimation using geometric priors in autonomous driving. *IEEE Robotics Automation Lett.* **7**, 11998–12005 (2022)
52. Xue, F., Zhuo, G., Huang, Z., Fu, W., Wu, Z., Ang, M.H.: Toward hierarchical self-supervised monocular absolute depth estimation for autonomous driving applications. In: 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 2330–2337 (2020). <https://doi.org/10.1109/IROS45743.2020.9340802>
53. Yin, Z., Shi, J.: Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In: CVPR (2018)
54. Zhan, H., Weerasekera, C.S., Bian, J.W., Reid, I.: Visual odometry revisited: What should be learnt? In: 2020 IEEE International Conference on Robotics and Automation (ICRA), pp. 4203–4210 (2020). <https://doi.org/10.1109/ICRA40945.2020.9197374>
55. Zhang, H., Li, Y., Cao, Y., Liu, Y., Shen, C., Yan, Y.: Exploiting temporal consistency for real-time video depth estimation. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 1725–1734 (2019). <https://doi.org/10.1109/ICCV.2019.00181>
56. Zhang, Y., et al.: Perceptual consistency in video segmentation. In: 2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pp. 2623–2632 (2022). <https://doi.org/10.1109/WACV51458.2022.00268>

57. Zhao, H., Gallo, O., Frosio, I., Kautz, J.: Loss functions for image restoration with neural networks. *IEEE Trans. Comput. Imaging* **3**(1), 47–57 (2017). <https://doi.org/10.1109/TCI.2016.2644865>
58. Zhao, W., Liu, S., Shu, Y., Liu, Y.J.: Towards better generalization: Joint depth-pose learning without posenet. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 9148–9158 (2020). <https://doi.org/10.1109/CVPR42600.2020.00917>
59. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6612–6619 (2017). <https://doi.org/10.1109/CVPR.2017.700>
60. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: 2017 IEEE International Conference on Computer Vision (ICCV), pp. 2242–2251 (2017). <https://doi.org/10.1109/ICCV.2017.244>
61. Zou, Y., Ji, P., Tran, Q.-H., Huang, J.-B., Chandraker, M.: Learning monocular visual odometry via self-supervised long-term modeling. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.-M. (eds.) *ECCV 2020*. LNCS, vol. 12359, pp. 710–727. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-58568-6_42
62. Zou, Y., Luo, Z., Huang, J.-B.: DF-Net: unsupervised joint learning of depth and flow using cross-task consistency. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) *ECCV 2018*. LNCS, vol. 11209, pp. 38–55. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-01228-1_3