# Consistency Guided Scene Flow Estimation

Yuhua Chen[1,2], Luc Van Gool[2], Cordelia Schmid[1], and Cristian Sminchisescu[1]

[1] Google Research
[2] ETH Zurich

**Abstract.** We present Consistency Guided Scene Flow Estimation (*CGSF*), a framework for joint estimation of 3D scene structure and motion from stereo videos. The model takes two temporal stereo pairs as input, and predicts disparity and scene flow. The model self-adapts at test time by iteratively refining its predictions. The refinement process is guided by a consistency loss, which combines stereo and temporal photo-consistency with a geometric term that couples the disparity and 3D motion. To handle the noise in the consistency loss, we further propose a learned, output refinement network, which takes the initial predictions, the loss, and the gradient as input, and efficiently predicts a correlated output update. We demonstrate with extensive experiments that the proposed model can reliably predict disparity and scene flow in many challenging scenarios, and achieves better generalization than the state-of-the-arts.

**Keywords:** scene flow, disparity estimation, stereo video, geometric constraints. self-supervised learning.

## 1 Introduction

Scene flow is the task of jointly estimating the 3D structure and motion of a scene [33]. This is critical for a wide range of downstream applications, such as robotics, autonomous driving and augmented reality. Typically, scene flow is estimated from consecutive stereo pairs, and requires simultaneously solving two sub-problems: stereo matching and optical flow. Though closely related, the two tasks cover different aspects of scene flow: stereo matching seeks to estimate disparity from a static stereo pair, whereas the objective of optical flow is to find the correspondence between temporally adjacent frames.

Powered by deep neural network, many end-to-end models have emerged for scene flow estimation [1,14] and its components [4,12,28], and showed promise on several benchmarks. However, training such models relies on the availability of labeled data. Due to difficulty in acquiring real-world data, synthetic data [21] is widely used as the major supervision source for deep models, with 3D structure or motion generated automatically using computer graphics techniques. Unfortunately, synthetic data still exhibits noticeable appearance dissimilarity and domain shift, which often prevents models trained this way from generalizing to the real-world, as shown in previous work [25,30]. This problem can be partially alleviated by finetuning on labeled real-world data. However, collecting ground-truth labels for scene flow can be extremely challenging, requiring the use of

costly sensors and additional manual intervention [23,24]. Therefore it may not always be feasible in many real-world scenarios, greatly limiting the applicability of deep models.

To address this issue, we present Consistency Guided Scene Flow (*CGSF*), a framework that additionally models output consistency. The framework begins with producing scene flow prediction using feedforward deep network. The predictions include disparity, optical flow and disparity change. These predictions are then coupled by a consistency loss which captures the photometric and geometric relations in scene flow.

Such consistency encodes the prior knowledge of scene flow, and is invariant to domain shift. Therefore, the consistency loss can be a powerful cue to guide the prediction for better generalization in unseen domains. However, the consistency loss is inherently noisy due to the complexity of natural data, such as non-Lambertian surfaces or occlusions. As a result, the provided guidance is not always precise, and can lead to undesired artifacts. To tackle this issue, we further propose a learned refinement network, which takes the initial predictions, the loss, and the gradient as input, and predicts an update to refine prediction recursively. The refinement network is implemented as a neural network, and thus can be trained jointly with the feedforward module.

Our *CGSF* model can be trained either using synthetic data, which can generalize well to real imagery, or trained on unlabeled data using the proposed consistency loss as self-supervision. In extensive experiments, we show our *CGSF* can reliably predict disparity and scene flow in many challenging scenarios, and the proposed learned refinement module can significantly improve the feedforward results, by maintaining the consistency in predictions. In particular, we demonstrate that the proposed model significantly improves the ability of generalization to unseen domains, which enables applications of real-world scene flow estimation.

## 2   Related Work

**Scene Flow Estimation.** As introduced by Vedula *et al.* [33], the process of scene flow estimation aims to recover the 3D structure and motion of a scene. The task has been formulated as a variational inference problem by multiple authors [2,11,34,36]. Recently, several deep learning based models have been built for the task. Ilg *et al.* [13] combined networks for stereo and optical flow using an additional network to predict disparity change. Ma *et al.* [20] stacked three networks for predicting disparity, flow and segmentation, then used a Gaussian-Newton solver to estimate per-object 3D motion. To leverage the correlation between tasks, a shared encoder architecture was proposed by Jiang *et al.* [14], where the encoder is shared among the tasks of disparity, optical flow, and segmentation. Similarly, Aleotti *et al.* [1] proposed a lightweight architecture to share information between tasks. Complementary to the previous work, our main objective in this work is to improve the generalization of the scene flow.

**Flow and Disparity Estimation.** As scene flow requires simultaneously reasoning about two tasks: disparity and flow estimation, it has been largely influenced by the techniques used in end-to-end stereo matching and optical flow. Recently significant progress has been made for both tasks.

For optical flow estimation, FlowNet [7] represents the first attempt in building end-to-end deep models. Then, FlowNet 2.0 [12] further improved the performance by dataset scheduling and refinement networks. A smaller network with a spatial pyramid is used in SpyNet [27] to handle large motions. Similar design is adopted in PWC-Net [28], which further includes cost volume processing at various pyramid levels. PWC-Net achieved highly competitive performance on several benchmark for optical flow estimation.

In stereo matching, Mayer *et al.* [21] introduced an architecture similar with FlowNet, for disparity estimation from rectified image pairs. 3D convolutions are used in GC-Net [15] are used to improve the matching accuracy. Similarly in PSM-Net [4], a pyramid pooling module is constructed to further exploit geometric context. More recently, GA-Net [38] integrates a semi-global aggregation layer and a local guided aggregation layer.

In this work, we build our feedforward module upon GA-Net [38] and PWC-Net [28] due to competitive performance.

**Adaptation.** Due to the difficulty of acquiring ground-truth labels in practice, the aforementioned models are often trained on synthetic data, and a performance drop is observed when applied in real-world [25,30].
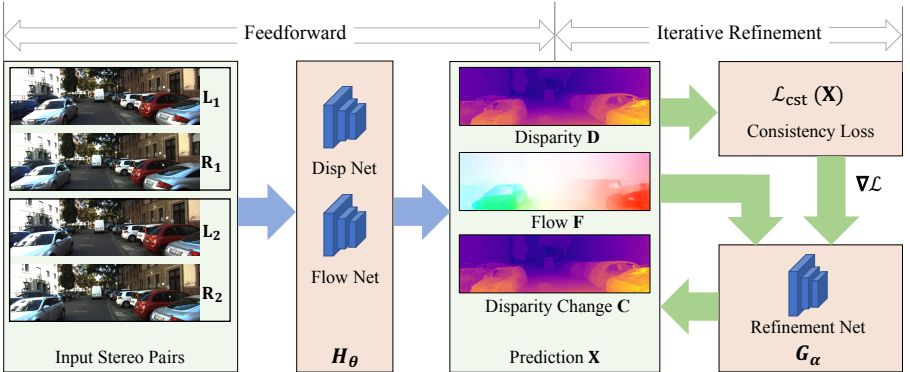
To address the issue, several techniques [10,26,30,31,32,39] have been proposed for adapting the model to new unlabelled domain. For the purpose, either offline adaptation [30] or online adaptation [32] is performed to update the model parameters. Our key difference in this work is that our refinement module operates on prediction directly instead of network parameters, and it thus has the advantage of being able to explicitly formulate a consistency which can guide the refinement process.

On the other hand, self-supervised learning has been used to learn from unlabelled data [3,8,9,17,35,37,41]. Typically, these methods requires collecting unlabelled data for offline training, which might not be always possible. Besides these trained model still suffer from the domain shift in the new untrained environments, and further online adaption is needed. Our model is oriented in similar principles, and our model can also be learned self-supervisedly in the offline training stage. But we additionally integrate the consistency within our refinement module, which makes our framework self-adapt to new environments without collecting unlabelled data beforehand or performing offline training.

**Online Iterative Refinement.** In terms of methodology, our work can be related to the line of work which aim to integrate iterative refinement into deep networks for various purposes, such as learning depth and pose estimation [6], 3D rigid motion estimation [19], monocular reconstruction of static scenes [29]. In contrast, our model consists of several novel geometric constraints for scene flow estimation, designed for self-supervision and overcoming domain shift.

# 3    Consistency Guided Scene Flow

## 3.1    Overview



**Fig. 1. Illustration of our proposed Consistency Guided Scene Flow (*CGSF*).** The self-supervised model can be trained either using labeled data, or unlabeled images using self-supervision. At run-time, it self-adapts by producing a set of outputs feedforward (networks $\mathbf{H}_\theta$, flow by blue arrows), then iteratively refining those estimates on two temporally adjacent stereo pairs, using a learned refinement scheme (network $\mathbf{G}_\alpha$, flow by green arrows). The refinement process is guided by the consistency loss, which captures intrinsic photometric and geometric constraints of scene flow.

In this section we present our Consistency Guided Scene Flow (*CGSF*) framework. An overview is given in fig. 1. Given as input a couple of stereo images $(\mathbf{L}_1, \mathbf{R}_1)$ and $(\mathbf{L}_2, \mathbf{R}_2)$ at subsequent timesteps, the model aims at predicting scene flow $\mathbf{X}$, *i.e.*, the 3D structure and motion of the environment. The 3D scene structure can be represented by the disparity in the first frame $\mathbf{D}_1$, as disparity is inversely proportional to depth. The 3D scene motion can be decomposed into the motion parallel with the image plane, represented by optical flow $\mathbf{F}_{1\rightarrow 2}$, and the motion perpendicular to the image plane, represented as disparity change $\mathbf{C}_{1\rightarrow 2}$.

At runtime, the scene flow estimation is formulated as an iterative optimization problem in this work. Our framework predicts scene flow in two stages: as feedforward pass and iterative output refinement. In the first stage, we rely on a feedforward module $\mathbf{H}_\theta$ to predict scene flow given input images. In the refinement stage, our geometric and photometric consistency losses (used during training) are further optimized with respect to the outputs using a refinement network $\mathbf{G}_\alpha$. Both the feedforward module $\mathbf{H}_\theta$ and the refinement network $\mathbf{G}_\alpha$ are implemented as neural networks with learnt parameters $\theta$ and $\alpha$. Therefore the two components can be trained jointly in an end-to-end process.

In the sequel, we introduce the feedforward module in §3.2. Then we formulate the geometric and photometric consistency terms for scene flow estimation

in §3.3. We present the learned refinement network in §3.4 and describe the training protocol in §3.5.

## 3.2   Feedforward Scene Flow Module

In the feedforward pass stage, the prediction of scene flow is obtained by passing the input images through the disparity and flow networks. The disparity network predicts the 3D structure given stereo pairs. Specifically, the module predicts $\mathbf{D}_1$ from $(\mathbf{L}_1, \mathbf{R}_1)$, and $\mathbf{D}_2$ from $(\mathbf{L}_2, \mathbf{R}_2)$. Given $(\mathbf{L}_1, \mathbf{L}_2)$ as input, the flow network predicts optical flow $\mathbf{F}_{1\rightarrow 2}$, and the disparity change $\mathbf{C}_{1\rightarrow 2}$ which encodes the motion component orthogonal to the image plane.

In this work, we build the feedforward module upon GA-Net [38] and PWC-Net [28], mainly due to the competitive performance and implementation availability. Here we briefly introduce the architectures.

**Disparity network**: We utilize GA-Net [38] to predict scene structure, *i.e.*, $\mathbf{D}_1$ and $\mathbf{D}_2$. The network extracts features from a stereo pair with a shared encoder. Then features are used to form a cost volume, which is fed into a cost aggregation block for regularization, refinement and disparity prediction.

**Flow network**: We modify PWC-Net [28] to predict the 3D motion, including optical flow $\mathbf{F}_{1\rightarrow 2}$ and disparity change $\mathbf{C}_{1\rightarrow 2}$. We follow the design in PWC-Net to construct a cost volume from the features of the first frame, and the warped features of the second frame. The cost volume is further processed to produce the motion prediction. We increase the output channel from 2 to 3 to encode 3D motion.

In principle, our framework is agnostic to specific choices of model design, thus other options of structure and motion networks are possible. This offers the flexibility of leveraging the latest advances in the field but still rely on the novel self-supervised geometric and photometric consistency losses, described next.

## 3.3   Scene Flow Consistency Loss

In this section we formulate the photometric and geometric constrains used in our proposed consistency loss.

**Stereo Photometric Consistency.** Photometric consistency between stereo pairs is widely used as a self-supervised loss in deep models [8,16]. The idea is to warp the source view via a differentiable bilinear sampling produced by the network prediction. The photometric difference between the warped source image and target image can be used as a self-supervision loss on the network prediction.

Given a stereo pair $(\mathbf{L}, \mathbf{R})$, we use the left view $\mathbf{L}$ as the target image, and the right view $\mathbf{R}$ as the source image. Considering a pixel $\mathbf{p}$ in the left view, its correspondence in the right view can be obtained by subtracting its disparity value from its x coordinate. With a slight abuse of notation, we denote it as $\mathbf{p} - \mathbf{D}(\mathbf{p})$. If $\mathbf{p}$ corresponds to a 3D scene point visible in both views, its correspondence

**Fig. 2. Producing occlusion map from disparity.** Given the original image (a) and its disparity prediction (c), a warped image (b) can be produced. Notice artifacts in the occluded areas. Given a scanline, as shown in green in (a), we plot in (e) the x-coordinates of the left and right correspondence. The occluded pixels have smaller disparity than the pixels with the same coordinates in the right image, therefore occluded pixels are marked in red. The inferred occlusion map is shown in (d).

should be visually consistent. Therefore the photometric stereo consistency can be written as

$$\mathcal{L}_{pd}(\mathbf{p}) = \mathbf{O}_d(\mathbf{p})||\mathbf{L}(\mathbf{p}) - \mathbf{R}(\mathbf{p} - \mathbf{D}(\mathbf{p}))|| \tag{1}$$

The stereo consistency applies to both $\mathbf{D}_1$ and $\mathbf{D}_2$. $\mathbf{O}_d$ is the occlusion map between left and right views, where $\mathbf{O}_d(\mathbf{p}) = 0$ when $\mathbf{p}$ is occluded and 1 otherwise. In this work we derive the occlusion map $\mathbf{O}_d$ from disparity $\mathbf{D}$ by means of a scanline algorithm, illustrated in fig. 2.

**Flow Photometric Consistency.** We also rely on photometric consistency for temporal frames. Given a pixel $\mathbf{p}$ in the first frame $\mathbf{L}_1$, its correspondence in the second frame $\mathbf{L}_2$ can be obtained by displacing $\mathbf{p}$ using optical flow. We can write the flow photometric consistency as

$$\mathcal{L}_{pf}(\mathbf{p}) = \mathbf{O}_f(\mathbf{p})||\mathbf{L}(\mathbf{p}) - \mathbf{L}(\mathbf{p} + \mathbf{F}_{1 \to 2}(\mathbf{p}))|| \tag{2}$$

$\mathbf{O}_f$ represents the occlusion map caused by optical flow. Similar with previous work [18,22], we infer the occlusion map by a forward-backward flow consistency check. In more detail, an additional backward flow is computed and then warped by the the forward flow $\mathbf{F}_{1 \to 2}$. We denote the computed warped backward flow map by $\hat{\mathbf{F}}_{2 \to 1}$. A forward-backward consistency check is performed to estimate the occlusion map

$$\mathbf{O}_f = [\![|\mathbf{F}_{1 \to 2} + \hat{\mathbf{F}}_{2 \to 1}|^2 < w_1(|\mathbf{F}_{1 \to 2}|^2 + |\hat{\mathbf{F}}_{2 \to 1}|^2) + w_2]\!] \tag{3}$$

where $[\![\cdot]\!]$ is the Iverson bracket, and we set $w_1 = 0.01, w_2 = 0.05$ to infer the occlusion mask.

**Disparity-Flow Consistency.** We formulate a new consistency term to connect disparity and 3D motion. Given a pixel $\mathbf{p}$ in the left view of the first frame $\mathbf{L}_1$, its correspondence in $\mathbf{L}_2$ is associated by optical flow $\mathbf{F}_{1\to2}(\mathbf{p})$ can be written as $\mathbf{p} + \mathbf{F}_{1\to2}(\mathbf{p})$. Thus its disparity value in the second frame should be $\mathbf{D}_2(\mathbf{p} + \mathbf{F}_{1\to2}(\mathbf{p}))$. On the other hand, $\mathbf{p}$'s disparity value in the second frame can also be obtained by adding the disparity change $\mathbf{C}_{1\to2}$ to the disparity of the first frame $\mathbf{D}_1$. Therefore the disparity flow consistency can be written as

$$\mathcal{L}_{df}(\mathbf{p}) = \mathbf{O}_f(\mathbf{p})||\mathbf{D}_1(\mathbf{p}) + \mathbf{C}_{1\to2}(\mathbf{p}) - \mathbf{D}_2(\mathbf{p} + \mathbf{F}_{1\to2}(\mathbf{p}))|| \qquad (4)$$

We use $\mathbf{O}_f$ to mask out occluded pixels as flow warping is used in the process.

**Smoothness Term.** To provide additional regularization of predictions, we use an edge-aware smoothness term used in [8]. The term is applied to the disparity map $\mathbf{D}_1, \mathbf{D}_2$, optical flow $\mathbf{F}_{1\to2}$, and disparity change $\mathbf{C}_{1\to2}$.

**Overall Consistency Loss.** By integrating all previously derived terms, the consistency loss can be written as

$$\mathcal{L}_{cst} = \mathcal{L}_{pd} + \mathcal{L}_{pf} + \mathcal{L}_{df} \qquad (5)$$

The consistency loss regularizes the prediction resulting from the individual structure and flow networks $\mathbf{X}$, and captures several intrinsic geometric and photometric constraints that predictions should satisfy.

### 3.4 Consistency Guided Refinement

In the refinement stage, the goal is to recursively improve scene flow prediction. We denote as $\mathbf{X}^t$ the scene flow prediction after $t$ steps of refinement. The initial prediction from the feedforward module can be denoted as $\mathbf{X}^0$.

As the consistency loss $\mathcal{L}_{cst}$ reflects the intrinsic structure of scene flow, it can be used as an indicator of how good the prediction is. Therefore we aim to leverage the signal from the consistency loss to guide the online refinement. To this end, we build a refinement network $\mathbf{G}_\alpha$ parameterized by $\alpha$, which takes as input the previous prediction $\mathbf{X}^t$, the consistency loss over the previous prediction $\mathcal{L}_{cst}(\mathbf{X}^t)$ and the gradient on predictions $\nabla_{\mathbf{X}}\mathcal{L}_{cst}(\mathbf{X}^t)$. The motivation of using the gradient is to provide a strong signal for the descent direction. The refinement network then predicts an update to refine prediction $\mathbf{X}^{t+1}$.

$$\mathbf{X}^{t+1} = \mathbf{X}^t + \mathbf{G}_\alpha(\mathbf{X}, \mathcal{L}_{cst}, \nabla_{\mathbf{X}}\mathcal{L}_{cst}) \qquad (6)$$

For implementation, we first concatenate all predictions, including $\mathbf{D}_1, \mathbf{D}_2, \mathbf{C}_{1\to2}$ each with size $W \times H \times 1$, and optical flow $\mathbf{F}_{1\to2}$ with size $W \times H \times 2$. The concatenated prediction map is of dimension $W \times H \times 5$, with $W$ and $H$ being the width and height of the image, respectively. The gradient is of the same dimensionality with the prediction, and the loss map is of 1 channel. Therefore,

all inputs can be concatenated as a $W \times H \times 11$ map. The benefit of using a concatenated map as input is that cross-channel correlations can be better captured. For instance, disparity information might be helpful to improve optical flow.

We implement refinement module as a small fully convolutional network(FCN) operating on the input map, to provide dense (per pixel) output for the updates. The advantage is that the FCN architecture can learn to leverage the local inter-pixel context, which is helpful to reduce noise in the consistency loss and propagate information to/from neighbouring pixels. We use a lightweight design of 3 convolution layers, each with 512 hidden units and kernel size of $3 \times 3$. ReLU is used between convolution layers as the activation function. The lightweighted design is motivated by the recursive use of the module.

### 3.5   Training

To learn the parameters of the feedforward scene flow $\mathbf{H}_\theta$ and refinement network $\mathbf{G}_\alpha$, the framework can be initialized based on pre-training using synthetic data with ground-truth labels (whenever available), and then/or further trained using the self-supervised loss $\mathcal{L}_{cst}$, in order to better generalize to real imagery. When the ground-truth label is available, a standard supervised loss writes

$$\mathcal{L}_{sup}(\mathbf{X}) = ||\mathbf{F}_{1\rightarrow2} - \mathbf{F}^*_{1\rightarrow2}|| + ||\mathbf{D}_1 - \mathbf{D}^*_1|| + ||\mathbf{C}_{1\rightarrow2} - \mathbf{C}^*_{1\rightarrow2}|| \tag{7}$$

where in the above equation, $*$ stands for ground-truth. In either case, the network parameters $\theta, \alpha$ can be learned jointly by minimizing the loss using stochastic gradient descent. The optimization objective writes

$$\arg\min_{\theta,\alpha} \sum_{t\in\{0,...,T\}} \mathcal{L}(\mathbf{X}^t) \tag{8}$$

where the loss $\mathcal{L}$ can be either $\mathcal{L}_{sup}$(when the label is available) or $\mathcal{L}_{cst}$(in self-supervised training). The loss is applied to the initial prediction, as well as to the prediction after each refinement step (i.e. we supervise the augmented step update at each iteration). Essentially after each step we minimize the loss w.r.t. $\boldsymbol{\alpha}$ so that the learnt descent direction produces as large of a loss decrease as possible, and aggregate in parallel over $T = 5$ refinement steps.

## 4   Experiments

In this section, we present experimental validation of the proposed *CGSF*. First, we test our model on synthetic data on FlyingThings3D in §4.1. Then in order to verify real-image generalization, we test the model on KITTI in §4.2. Then we provide in §4.3 ablation studies to justify our design choices. Finally, in §4.4 we show the performance of *CGSF* on daily scenes captured by a stereo camera, thus to demonstrate the generalization to unseen environments.

| | Disparity EPE | Flow EPE | Disparity Change EPE * |
|---|---|---|---|
| DWARF [1] | 2.00 | 6.52 | 2.23 |
| Feedforward | 0.83 | 7.02 | 2.02 |
| $CGSF$(ours) | **0.79** | **5.98** | **1.33** |

**Table 1. Scene flow results on FlyingThings3D test set**. Results are evaluated as end point error (EPE). All models are trained on FlyingThings3D training set. *Disparity change EPE represents the error of the motion orthogonal to the image plane.

## 4.1  Evaluation on Synthetic Data

First we evaluate the model on synthetic data. We use the FlyingThings3D [21] which is a synthetic dataset for scene flow estimation, containing $22,390$ stereo training pairs. Dense ground-truth labels are available in all images for disparity, optical flow and disparity change. We use the ground-truth to supervise our model, including the feedforward module and the refinement network. In the refinement module, we use a set of trade-off weights to balance different consistency terms. We use $\omega_{pd}, \omega_{pf}, \omega_{df}$ to denote respectively the weight of stereo photometric consistency $\mathcal{L}_{pd}$, flow disparity consistency $\mathcal{L}_{pf}$, and disparity disparity-flow consistency $\mathcal{L}_{df}$. The weight of smoothness term is denoted as $\omega_s$. We set $\omega_{pd} = 1, \omega_{pf} = 1, \omega_{df} = 1, \omega_s = 0.1$. A stage-wise training strategy is employed to improve the training stability. In more details, we first train the disparity network and the flow network separately. For this part we follow [38] and [28] to set the hyperparameters. These are then used as initialization for further joint fine-tuning. Specifically, we randomly initialize the weights in refinement module, and jointly fine-tune the whole network including the pre-trained feedforward networks. To fine-tune the network, we use the Adam optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ for the training. The learning rate is set to $10^{-4}$. The network is trained on FlyingThings 3D for 10 epochs with a batch size of 4.

The models are evaluated on the test set of FlyingThings3D, which contains $4,370$ images. As scene flow model requires 4 images as input, thus we drop the last frame of every sequence, resulting in a test set with $3,933$ images. The performance is measured by end point error(EPE) in three outputs: optical flow, disparity, and disparity change. In the refinement module, we use a set of trade-off weights to balance different consistency terms. We use $\omega_{pd}, \omega_{pf}, \omega_{df}$ to denote respectively the weight of stereo photometric consistency $\mathcal{L}_{pd}$, flow disparity consistency $\mathcal{L}_{pf}$, and disparity disparity-flow consistency $\mathcal{L}_{df}$. The weight of smoothness term is denoted as $\omega_s$. We set $\omega_{pd} = 1, \omega_{pf} = 1, \omega_{df} = 1, \omega_s = 0.1$. We use the results produced by the feedforward module without a refinement network as baseline. To facilitate the comparison with the state-of-the-art, we test and compare with a recent model DWARF [1] in the same experimental setting. As shown in table 1, the refinement model improves over the feedforward baseline in all three metrics, suggesting the effectiveness of our consistency

| Method | Training Data | D1 | D2 | F1 | SF |
|---|---|---|---|---|---|
| Zhou *et al* [40] | K(u) | 9.41 | - | - | - |
| Godard *et al* [8] | K(u) | 9.19 | - | - | - |
| BridgeDepthFlow [17] | K(u) | 9.21 | - | 27.02 | - |
| UnOS [35] | K(u) | 5.90 | - | 20.12 | - |
| MADNet [32] | F(s)+ K(u) | 8.41 | - | - | - |
| DWARF [1] | F(s) | 11.60 | 29.64 | 38.32 | 45.58 |
| Feedforward | F(s) | 11.53 | 27.83 | 33.09 | 43.55 |
| *CGSF*(ours) | F(s) | 7.10 | 19.68 | 27.35 | 35.40 |
| *CGSF*(ours) | K(u) | 6.65 | 17.69 | 23.05 | 31.52 |
| *CGSF*(ours) | F(s)+ K(u) | 5.75 | 15.53 | 20.45 | 28.14 |

**Table 2. Evaluation on KITTI 2015 scene flow dataset.** The percentage of outliers is used as evaluation metric. A pixel is considered as correct if the prediction end-point error is smaller than $3px$ or 5%. For training data, 'K(u)' stands for unsupervised training on KITTI, and 'F(s)' represents supervised training on FlyingThings3D.
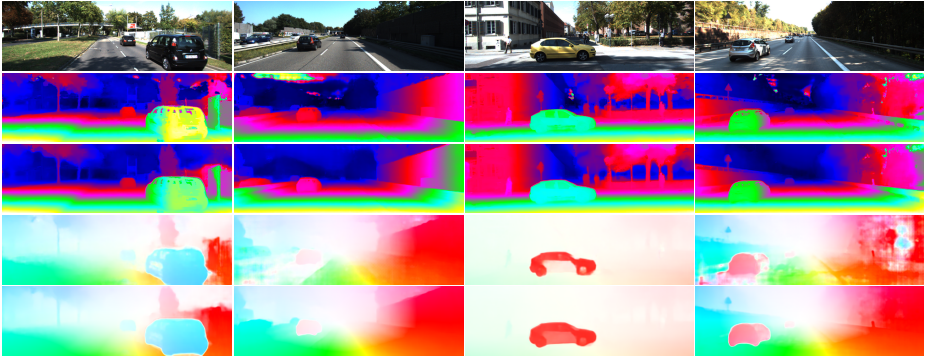
guided refinement. Our model also outperforms DWARF [1] by a large margin. Again, this shows the benefit of geometrically modelling the scene flow consistency among outputs.

### 4.2   Generalization to Real Imagery

To verify the performance on real imagery where no labeled data is available, we proceed the evaluation with KITTI scene flow dataset [23,24]. KITTI scene flow contains 200 training images with sparse ground-truth disparity and flow acquired by lidar sensor. The ground-truth is only used in evaluation. To comply with the KITTI scene flow benchmark, the percentage of outliers in the following 4 categories is used as the evaluation metrics, D1: outliers in the disparity of first frame (*i.e.*, $\mathbf{D}_1$), D2: outliers in second frame warped to the first frame (*i.e.*, $\mathbf{D}_1 + \mathbf{C}_{1\rightarrow 2}$), F1: outliers in optical flow (*i.e.*, $\mathbf{F}_{1\rightarrow 2}$), SF: outliers in either D1,D2 or F1. A pixel is considered correct if the prediction end-point error is smaller than $3px$ or 5%.

As discussed in §3.3, our *CGSF* model can be trained with the ground-truth labels from synthetic data, or using self-supervision. We evaluate both strategies. For training using ground-truth on synthetic data, we use FlyingThings3D as in the previous experiment. For self-supervised training on KITTI, we use the KITTI raw dataset as training set, with all test scenes excluded.

We compare with the supervised feedforward network without refinement, as a baseline. Additionally, we compare with several recent competing methods, including self-supervised stereo models [8,40], joint self-supervised stereo and flow models [17,35], scene flow model supervised on synthetic data [1], and an adaptation model for stereo estimation [32]. For [32], we initialize it with the pre-trained weights from synthetic data, we conduct unsupervised adaptations on KITTI raw dataset (the same set used in our unsupervised training).

**Fig. 3. Disparity estimation results on KITTI.** From top to bottom: input image, disparity from feedforward module, disparity after refinement, flow from feedforward module, flow after refinement. Note that our refinement module improves the quality of both flow and disparity estimation.

The results are summarized in table 2. With the feedforward module only, the baseline achieves an error rate of 43.55%. The error can be significantly reduced to 35.40% by the proposed refinement network. Notably, the performance gain is much larger compared to the one in FlyingThings3D, which demonstrates that consistency plays a central and positive role in cross-domain scenarios. This indicates that consistency can be used as an effective model to overcome domain gaps and improve generalization. Some qualitative results are displayed in fig. 3, where a considerable improvement in both disparity and flow estimation is observed, by using the proposed consistency guided refinement. This again confirms the benefit of explicitly modelling consistency in the refinement process.

Whenever collecting an unlabeled training set beforehand is possible, the model can also be trained under self-supervision, results an error rate of 31.52%, which demonstrates the effectiveness of our proposed consistency loss as a self-supervised loss, and also the good compatibility of our model in terms of supervision. The performance can be further improved by combining pre-training on synthetic data, and using self-supervision for finetuning. This combination results in the lowest scene flow error rate of 28.14%.
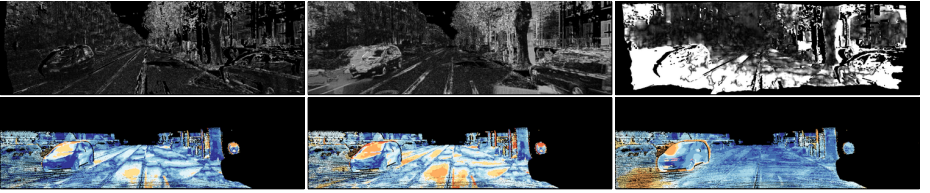
### 4.3 Ablation Studies

We provide further experimental analysis on ablated versions of our model, in order to justify the design choices. To ensure a fair comparison, all models used in this section are trained on FlyingThings3D, and results are reported on KITTI2015 without further self-supervised finetuning.

**Consistency in Refinement Network** In the refinement network, the consistency loss is used as a guidance to refine the network outputs. The consistency

| Refinement | $\mathcal{L}_{pd}$ | $\mathcal{L}_{pf}$ | $\mathcal{L}_{df}$ | D1 | D2 | F1 | SF |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | | | | 11.29 | 26.51 | 32.71 | 43.01 |
| ✓ | | | | 10.40 | 25.15 | 32.05 | 42.43 |
| ✓ | ✓ | | | 7.95 | 22.48 | 31.38 | 39.80 |
| ✓ | | ✓ | | 9.84 | 23.92 | 29.06 | 39.64 |
| ✓ | | | ✓ | 9.11 | 21.14 | 30.95 | 39.59 |
| ✓ | ✓ | ✓ | | 7.60 | 22.73 | 29.07 | 38.23 |
| ✓ | ✓ | ✓ | ✓ | 7.10 | 19.68 | 27.35 | 35.40 |

**Table 3. Ablation study on consistency terms.** Different consistency terms are used in the refinement module. All models are trained on FlyingThings3D and evaluated on KITTI.



**Fig. 4. Illustration of different consistency terms.** From left to right we show results produced using $\mathcal{L}_{pd}, \mathcal{L}_{pf}, \mathcal{L}_{df}$. The first row presents the pixel-wise loss map on the feedforward prediction, where brighter is larger loss. The second row shows the error map for disparity estimation, when compared with the ground-truth label. Red is higher error, and blue lower error.

loss mainly consists of three terms: stereo photometric loss ($\mathcal{L}_{pd}$), flow photometric loss ($\mathcal{L}_{pf}$) and disparity flow consistency loss ($\mathcal{L}_{df}$). To understand the role of each loss, we conduct an ablation study where we use different combinations in the refinement network.

We summarize the results in table 3. When no consistency is used, the refinement network only takes the prediction $\mathbf{X}$ as input. This baseline results in a scene flow error of of 42.43, slightly better than the feedforward result which is 43.01. This suggests that only $\mathbf{X}$ is insufficient in effectively guiding refinement. The error can be reduced by additionally including consistency terms, validating the efficacy of our guided refinement strategy. The contributions of each loss term to the scene flow error are similar. However, each loss term exhibits different improvements over the three sub-measurements. Notably, $\mathcal{L}_{pd}$ improves $D1$ measure the most, $\mathcal{L}_{pf}$ improves $F1$ the most, while $\mathcal{L}_{df}$ reduces $D2$ error, which is related to both disparity and optical flow. The results are understandable, as each loss intrinsically constraints different outputs. For example, the stereo photometric loss $\mathcal{L}_{pd}$ links stronger to disparity, whereas the flow photometric loss $\mathcal{L}_{pf}$ better constrains flow.

To further understand the impact of different losses, we visualize refinement results by using only $\mathcal{L}_{pd}$, $\mathcal{L}_{pf}$ or $\mathcal{L}_{df}$ in fig. 4. We observe that the photometric

| $\mathbf{X}$ | $\mathcal{L}$ | $\nabla\mathcal{L}$ | D1 | D2 | F1 | SF |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| ✓ | | | 10.40 | 25.15 | 32.05 | 42.43 |
| ✓ | ✓ | | 9.77 | 22.80 | 30.95 | 40.47 |
| ✓ | | ✓ | 7.44 | 20.01 | 28.20 | 36.31 |
| ✓ | ✓ | ✓ | 7.10 | 19.68 | 27.35 | 35.40 |

**Table 4. Ablation study on the input to the refinement network.** All models are trained on FlyingThings3D, and evaluated on KITTI.

loss (the left, and the middle) is sparse and cannot provide sufficient signal in the textureless regions. As a result, the error in flat regions (such as road) remains large after refinement. On the other hand, the disparity flow consistency term $\mathcal{L}_{df}$ provides a denser signal, especially in textureless regions, and can reduce errors there. However, it still produces errors on the motion boundaries and occluded parts, due to output inconsistency in such regions.

**Input to Refinement Network** In the refinement network, consistency is used to guide the refinement of scene flow prediction, in which we use the current scene flow prediction $\mathbf{X}^t$, the consistency computed on the scene flow $\mathcal{L}_{cst}(\mathbf{X})$, and its gradient as input $\mathcal{L}_{cst}(\mathbf{X})$. In this study, we investigate different inputs and the influence on final performance. As shown in table 4, with only the variable as input, the refinement network achieved 42.43. By additionally including the loss and gradient as inputs, a further decrease of 0.96 and 6.12 are observed, respectively. The results suggest that the gradient is a stronger signal for refinement, compared to the loss. Nevertheless, combining the two results in the best performance in scene flow, at 35.40 error rate.
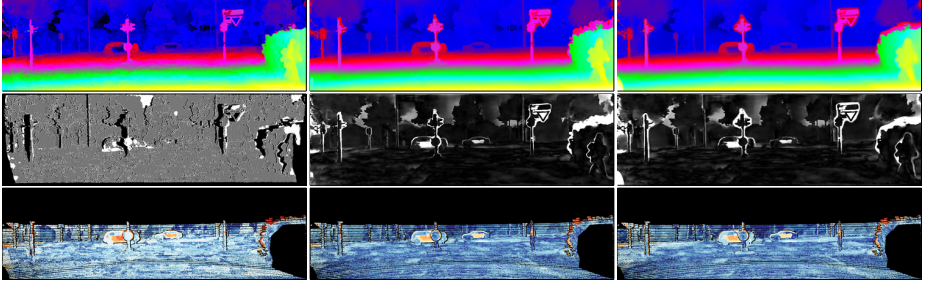
**Online Refinement Strategies** In this experiment we compare our learned refinement module against alternative design choices. In this work we use a refinement network to predict an update from the loss signal, and to refine the scene flow prediction without updating the model parameters. Alternatively, one can perform online finetuning to update the model parameters, based on the same consistency loss. We refer to this approach as parameter finetuning. Another choice is to change the scene flow prediction directly, without passing through the network, which we refer to as output finetuning [5].

We test the two alternatives with the same consistency loss. The results are summarized in table 5. As can be seen from the table, our learned refinement achieves similar performance with parameter finetuning. However, as multiple iterations of forward and backward passes are needed for finetuning network parameters, this approach results in a large run time of 120 seconds per image. Compared to parameter finetuning, our method is much faster.

To further understand the difference of the three refinement strategies, we visualize the disparity outputs in fig. 5. We can observe from the figure, that

| Refinement Method | D1 | D2 | F1 | SF | Time |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Outputs | 8.53 | 22.82 | 30.61 | 39.95 | 0.5s |
| Parameters | 6.88 | 19.47 | 27.69 | 35.49 | 120s |
| Learned (ours) | 7.10 | 19.68 | 27.35 | 35.40 | 0.8s |

**Table 5. Ablation study on online refinement strategies.** We use different online refinement strategies to refine scene flow estimation. All models are trained on FlyingThings3D and results reported on KITTI.
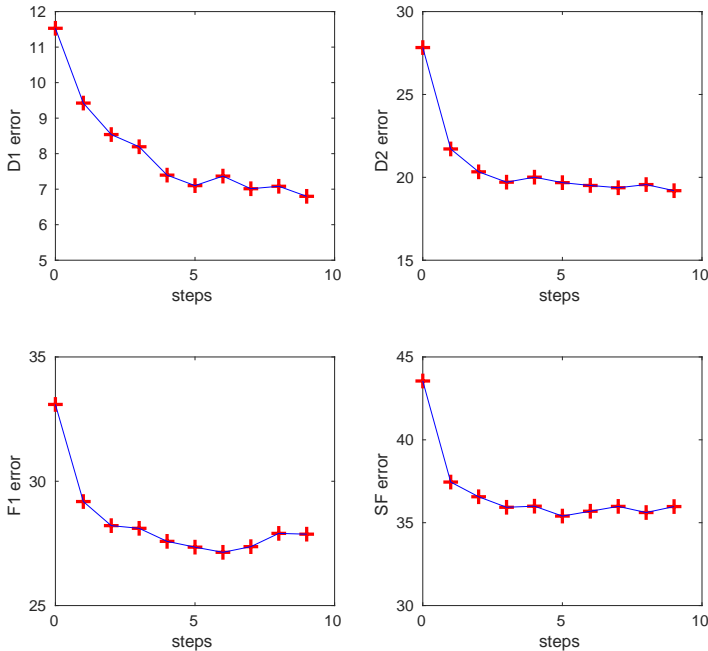


**Fig. 5. Qualitative results for different online refinement strategies.** From left to right we show the results refined with output finetuning, parameter finetuning and our proposed learning-based refinement. The first row illustrates disparity estimation, the second the accumulated updates, where brighter means larger change in disparity. The last row shows the error map evaluated against ground-truth. Red represents higher error, and blue lower error.

output finetuning is very noisy and sparse – this is due to properties of each pixel being refined independently. As shown previously, the loss signal is sparse and noisy, and is often invalid in occluded regions, which results the noisy updates. In contrast, our learned refinement can take advantage of such sparse signal, and learn to propagate the guidance signal by leveraging context information. As a result, a much smoother refinement (similar to the one produced by parameter finetuning, but two orders of magnitude faster) is produced by our learned refinement module.

**Effect of Refinement Module** The consistency-guided refinement module is used to iteratively refine the scene flow estimation. To gain further understanding of the refinement process, we provide an analysis on the intermediate results, *i.e.*, scene flow estimation after each refinement step.

In more details, we evaluate the scene flow estimation results after each refinement step, using the four evaluation metrics: D1, D2, F1, SF, which represent the percentage of outliers. In the main paper, a total of 5 refinement steps are performed. Here we also report the performance of applying further refinement, up to 10 steps. We plot the results at different steps in fig 6. We observe that
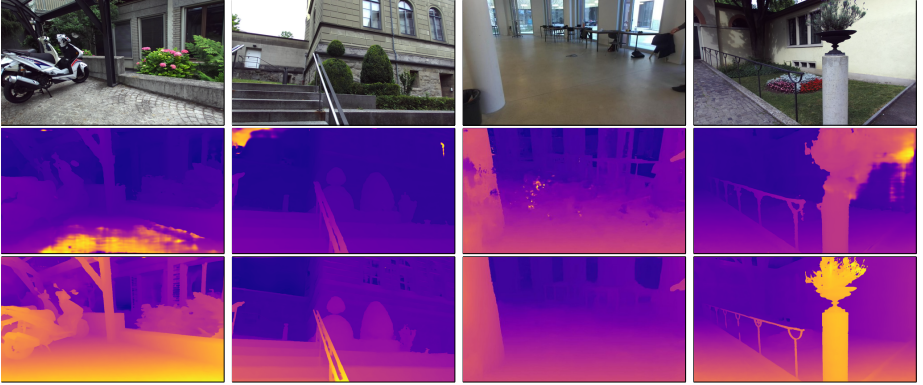
**Fig. 6. The effect of refinement steps on scene flow performance.** We plot the error using D1, D2, F1 and SF measure in the above four figures, respectively. x-axis is the number of refinement steps, with 0 being the feedforward results(without the refinement module). y-axis is the percentage of outliers. A pixel is considered correct if the prediction end-point error is smaller than $3px$ or 5%. For SF, a correct pixel needs to be correct for D1, D2 and F1. The model is trained on FlyingThings3D, and test on KITTI scene flow dataset.

the error decreases notably in all four measurements, which shows that the refinement module is useful for improving scene flow accuracy. The improvement is more significant in the first 5 steps, then the performance saturates with additional steps. Further refinement leads to an increase in the overall running time.

## 4.4   Performance in Daily Scenes

To test the generalization of our method in other environments, we apply the model to video captured by a ZED stereo camera. We compare the results from the feedforward module. As shown in fig. 7, our consistency guided refinement can notably improve the performance of disparity estimation in such case. This

**Fig. 7. Disparity estimation results in daily scenes.** Top: input left image, Middle: feedforward results, Bottom: results with consistency guided refinement. Note the refinement can significantly reduce drastically incorrect predictions.

demonstrates that consistency in scene flow is a powerful tool to overcome domain shift, and guide the model to generalize to unseen scenarios.

## 5  Conclusions

We have presented Consistency Guided Scene Flow (*CGSF*), a framework for joint estimation of disparity and scene flow from stereo video. The consistency loss combines stereo and temporal photo-consistency with a novel geometric consistency which couples the depth and flow variables. Besides the usual online parameter and output finetuning strategies typical of self-supervised test-time domain adaptation, we introduce new learned output finetuning descent models that produce good quality results, on par with parameter finetuning, but two orders of magnitude faster. Extensive experimental validation on benchmarks shows that *CGSF* can reliably predict disparity and scene flow, with good generalization in cross-domain settings.

# References

1. Aleotti, F., Poggi, M., Tosi, F., Mattoccia, S.: Learning end-to-end scene flow by distilling single tasks knowledge. In: AAAI (2020)
2. Basha, T., Moses, Y., Kiryati, N.: Multi-view scene flow estimation: A view centered variational approach. International journal of computer vision **101**(1), 6–21 (2013)
3. Casser, V., Pirk, S., Mahjourian, R., Angelova, A.: Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In: AAAI (2019)
4. Chang, J., Chen, Y.: Pyramid stereo matching network. In: CVPR (2018)
5. Chen, Y., Schmid, C., Sminchisescu, C.: Self-supervised learning with geometric constraints in monocular video: Connecting flow, depth, and camera. In: ICCV (2019)
6. Clark, R., Bloesch, M., Czarnowski, J., Leutenegger, S., Davison, A.J.: Ls-net: Learning to solve nonlinear least squares for monocular stereo. arXiv preprint arXiv:1809.02966 (2018)
7. Dosovitskiy, A., Fischer, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., Van Der Smagt, P., Cremers, D., Brox, T.: FlowNet: Learning optical flow with convolutional networks. In: ICCV (2015)
8. Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: CVPR (2017)
9. Godard, C., Mac Aodha, O., Firman, M., Brostow, G.J.: Digging into self-supervised monocular depth estimation. In: ICCV (2019)
10. Guo, X., Li, H., Yi, S., Ren, J., Wang, X.: Learning monocular depth by distilling cross-domain stereo networks. In: ECCV (2018)
11. Huguet, F., Devernay, F.: A variational method for scene flow estimation from stereo sequences. In: ICCV (2007)
12. Ilg, E., Mayer, N., Saikia, T., Keuper, M., Dosovitskiy, A., Brox, T.: FlowNet 2.0: Evolution of optical flow estimation with deep networks. In: CVPR (2017)
13. Ilg, E., Saikia, T., Keuper, M., Brox, T.: Occlusions, motion and depth boundaries with a generic network for disparity, optical flow or scene flow estimation. In: ECCV (2018)
14. Jiang, H., Sun, D., Jampani, V., Lv, Z., Learned-Miller, E., Kautz, J.: Sense: A shared encoder network for scene-flow estimation. In: ICCV (2019)
15. Kendall, A., Martirosyan, H., Dasgupta, S., Henry, P., Kennedy, R., Bachrach, A., Bry, A.: End-to-end learning of geometry and context for deep stereo regression. In: ICCV (2017)
16. Kuznietsov, Y., Stuckler, J., Leibe, B.: Semi-supervised deep learning for monocular depth map prediction. In: CVPR (2017)
17. Lai, H.Y., Tsai, Y.H., Chiu, W.C.: Bridging stereo matching and optical flow via spatiotemporal correspondence. In: CVPR (2019)
18. Liu, P., Lyu, M., King, I., Xu, J.: Selflow: Self-supervised learning of optical flow. In: CVPR (2019)
19. Lv, Z., Dellaert, F., Rehg, J.M., Geiger, A.: Taking a deeper look at the inverse compositional algorithm. In: CVPR (2019)
20. Ma, W.C., Wang, S., Hu, R., Xiong, Y., Urtasun, R.: Deep rigid instance scene flow. In: CVPR (2019)
21. Mayer, N., Ilg, E., Häusser, P., Fischer, P., Cremers, D., Dosovitskiy, A., Brox, T.: A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In: CVPR (2016)

22. Meister, S., Hur, J., Roth, S.: UnFlow: Unsupervised learning of optical flow with a bidirectional census loss. In: AAAI (2018)
23. Menze, M., Heipke, C., Geiger, A.: Joint 3d estimation of vehicles and scene flow. In: ISPRS Workshop on Image Sequence Analysis (ISA) (2015)
24. Menze, M., Heipke, C., Geiger, A.: Object scene flow. ISPRS Journal of Photogrammetry and Remote Sensing (JPRS) (2018)
25. Pang, J., Sun, W., Yang, C., Ren, J., Xiao, R., Zeng, J., Lin, L.: Zoom and learn: Generalizing deep stereo matching to novel domains. In: CVPR (2018)
26. Poggi, M., Pallotti, D., Tosi, F., Mattoccia, S.: Guided stereo matching. In: CVPR (2019)
27. Ranjan, A., Black, M.J.: Optical flow estimation using a spatial pyramid network. In: CVPR (2017)
28. Sun, D., Yang, X., Liu, M.Y., Kautz, J.: PWC-net: CNNs for optical flow using pyramid, warping, and cost volume. In: CVPR (2018)
29. Tang, C., Tan, P.: Ba-net: Dense bundle adjustment network. arXiv preprint arXiv:1806.04807 (2018)
30. Tonioni, A., Poggi, M., Mattoccia, S., Di Stefano, L.: Unsupervised adaptation for deep stereo. In: ICCV (2017)
31. Tonioni, A., Rahnama, O., Joy, T., Stefano, L.D., Ajanthan, T., Torr, P.H.: Learning to adapt for stereo. In: CVPR (2019)
32. Tonioni, A., Tosi, F., Poggi, M., Mattoccia, S., Stefano, L.D.: Real-time self-adaptive deep stereo. In: CVPR (2019)
33. Vedula, S., Baker, S., Rander, P., Collins, R., Kanade, T.: Three-dimensional scene flow. In: ICCV (1999)
34. Vogel, C., Schindler, K., Roth, S.: Piecewise rigid scene flow. In: ICCV (2013)
35. Wang, Y., Wang, P., Yang, Z., Luo, C., Yang, Y., Xu, W.: Unos: Unified unsupervised optical-flow and stereo-depth estimation by watching videos. In: CVPR (2019)
36. Wedel, A., Rabe, C., Vaudrey, T., Brox, T., Franke, U., Cremers, D.: Efficient dense scene flow from sparse or dense stereo data. In: ECCV (2008)
37. Yin, Z., Shi, J.: GeoNet: Unsupervised learning of dense depth, optical flow and camera pose. In: CVPR (2018)
38. Zhang, F., Prisacariu, V., Yang, R., Torr, P.H.: Ga-net: Guided aggregation net for end-to-end stereo matching. In: CVPR (2019)
39. Zhong, Y., Li, H., Dai, Y.: Open-world stereo video matching with deep rnn. In: ECCV (2018)
40. Zhou, C., Zhang, H., Shen, X., Jia, J.: Unsupervised learning of stereo matching. In: ICCV (2017)
41. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: CVPR (2017)