

Distilled Semantics for Comprehensive Scene Understanding from Videos

Fabio Tosi* Filippo Aleotti* Pierluigi Zama Ramirez*
 Matteo Poggi Samuele Salti Luigi Di Stefano Stefano Mattoccia
 Department of Computer Science and Engineering (DISI)
 University of Bologna, Italy

*{fabio.tosi5, filippo.aleotti2, pierluigi.zama}@unibo.it

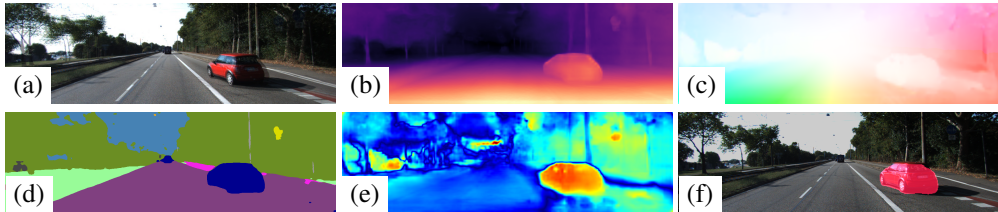


Figure 1. Given an input monocular video (a), our network can provide the following outputs in real-time: depth (b), optical flow (c), semantic labels (d), per-pixel motion probabilities (e), motion mask (f).

Abstract

Whole understanding of the surroundings is paramount to autonomous systems. Recent works have shown that deep neural networks can learn geometry (depth) and motion (optical flow) from a monocular video without any explicit supervision from ground truth annotations, particularly hard to source for these two tasks. In this paper, we take an additional step toward holistic scene understanding with monocular cameras by learning depth and motion alongside with semantics, with supervision for the latter provided by a pre-trained network distilling proxy ground truth images. We address the three tasks jointly by a) a novel training protocol based on knowledge distillation and self-supervision and b) a compact network architecture which enables efficient scene understanding on both power hungry GPUs and low-power embedded platforms. We thoroughly assess the performance of our framework and show that it yields state-of-the-art results for monocular depth estimation, optical flow and motion segmentation.

1. Introduction

What information would an autonomous agent be keen to gather from its sensory sub-system to tackle tasks like navigation and interaction with the explored environment? It would need to be informed about the geometry of the surroundings and the type of objects therein, and likely better

know which of the latter are actually moving and how they do so. What if all such cues may be provided by as simple a sensor as a single RGB camera?

Nowadays, deep learning is advancing the state-of-the-art in classical computer vision problems at such a quick pace that single-view holistic scene understanding seems to be no longer out-of-reach. Indeed, highly challenging problems such as monocular depth estimation and optical flow can nowadays be addressed successfully by deep neural networks, often through unified architectures [88, 3, 96]. Self-supervised learning techniques have yielded further major achievements [95, 58] by enabling effective training of deep networks without annotated images. In fact, labels are hard to source for depth estimation due to the need of active sensors and manual filtering, and are even more cumbersome in the case of optical flow. Concurrently, semi-supervised approaches [90, 16] proved how a few semantically labelled images can improve monocular depth estimation significantly. These works have also highlighted how, while producing per-pixel class labels is tedious yet feasible for a human annotator, manually endowing images with depth and optical flow ground-truths is prohibitive.

In this paper, we propose the first-ever framework for comprehensive scene understanding from monocular videos. As highlighted in Figure 1, our multi-stage network architecture, named Ω Net, can predict depth, semantics, optical flow, per-pixel motion probabilities and motion masks. This comes alongside with estimating the pose between adjacent frames for an uncalibrated camera, whose intrinsic parameters are also estimated. Our training methodology

*Joint first authorship.

leverages on self-supervision, knowledge distillation and multi-task learning. In particular, peculiar to our proposal and key to performance is distillation of proxy semantic labels gathered from a state-of-the-art pre-trained model [52] within a self-supervised and multi-task learning procedure addressing depth, optical flow and motion segmentation. Our training procedure also features a novel and effective self-distillation schedule for optical flow mostly aimed at handling occlusions and relying on tight integration of rigid flow, motion probabilities and semantics. Moreover, Ω Net is lightweight, counting less than 8.5M parameters, and fast, as it can run at nearly 60 FPS and 5 FPS on an NVIDIA Titan Xp and a Jetson TX2, respectively. As vouched by thorough experiments, the main contributions of our work can be summarized as follows:

- The first real-time network for joint prediction of depth, optical flow, semantics and motion segmentation from monocular videos
- A novel training protocol relying on proxy semantics and self-distillation to effectively address the self-supervised multi-task learning problem
- State-of-the-art self-supervised monocular depth estimation, largely improving accuracy at long distances
- State-of-the-art optical flow estimation among monocular multi-task frameworks, thanks to our novel occlusion-aware and semantically guided training paradigm
- State-of-the-art motion segmentation by joint reasoning about optical-flow and semantics

2. Related Work

We review previous works relevant to our proposal.

Monocular depth estimation. At first, depth estimation was tackled as a supervised [24, 49] or semi-supervised task [48]. Nonetheless, self-supervision from image reconstruction is now becoming the preferred paradigm to avoid hard to source labels. Stereo pairs [25, 28] can provide such supervision and enable scale recovery, with further improvements achievable by leveraging on trinocular assumptions [64], proxy labels from SGM [76, 80] or guidance from visual odometry [2]. Monocular videos [95] are a more flexible alternative, although they do not allow for scale recovery and mandate learning camera pose alongside with depth. Recent developments of this paradigm deal with differentiable direct visual odometry [77] or ICP [57] and normal consistency [87]. Similarly to our work, [88, 96, 17, 3, 86, 56] model rigid and non-rigid components using the projected depth, relative camera transformations, and optical flow to handle independent motions, which can also be estimated independently in the 3D space [9, 83]. In [30], the authors show how to learn camera intrinsics together with depth and egomotion to enable training on any unconstrained video. In [29, 94, 6], reasoned design choices such as a minimum reprojection loss between frames, self-

assembled attention modules and auto-mask strategies to handle static camera or dynamic objects proved to be very effective. Supervision from stereo and video have also been combined [91, 29], possibly improved by means of proxy supervision from stereo direct sparse odometry [84]. Uncertainty modeling for self-supervised monocular depth estimation has been studied in [63]. Finally, lightweight networks aimed at real-time performance on low-power systems have been proposed within self-supervised [62, 61] as well as supervised [81] learning paradigms.

Semantic segmentation. Nowadays, fully convolutional neural networks [55] are the standard approach for semantic segmentation. Within this framework, multi-scale context modules and proper architectural choices are crucial to performance. The former rely on spatial pyramid pooling [31, 93] and atrous convolutions [14, 13, 15]. As for the latter, popular backbones [47, 74, 32] have been improved by more recent designs [34, 18]. While for years the encoder-decoder architecture has been the most popular choice [70, 4], recent trends in Auto Machine Learning (AutoML) [52, 12] leverage on architectural search to achieve state-of-the-art accuracy. However, these latter have huge computational requirements. An alternative research path deals with real-time semantic segmentation networks. In this space, [60] deploys a compact and efficient network architecture, [89] proposes a two paths network to attain fast inferences while capturing high resolution details. DABNet [50] finds an effective combinations of depth-wise separable filters and atrous-convolutions to reach a good trade-off between efficiency and accuracy. [51] employs cascaded sub-stages to refine results while FCHardNet [11] leverages on a new harmonic densely connected pattern to maximize the inference performance of larger networks.

Optical flow estimation. The optical flow problem concerns estimation of the apparent displacement of pixels in consecutive frames, and it is useful in various applications such as, *e.g.*, video editing [10, 43] and object tracking [82]. Initially introduced by Horn and Schunck [33], this problem has traditionally been tackled by variational approaches [8, 7, 69]. More recently, Dosovitskiy *et al.* [21] showed the supremacy of deep learning strategies also in this field. Then, other works improved accuracy by stacking more networks [38] or exploiting traditional pyramidal [65, 75, 35] and multi-frame fusion [67] approaches. Unfortunately, obtaining even sparse labels for optical flow is extremely challenging, which renders self-supervision from images highly desirable. For this reason, an increasing number of methods propose to use image reconstruction and spatial smoothness [41, 68, 73] as main signals to guide the training, paying particular attention to occluded regions [58, 85, 53, 54, 40, 37].

Semantic segmentation and depth estimation. Monocular depth estimation is tightly connected to the

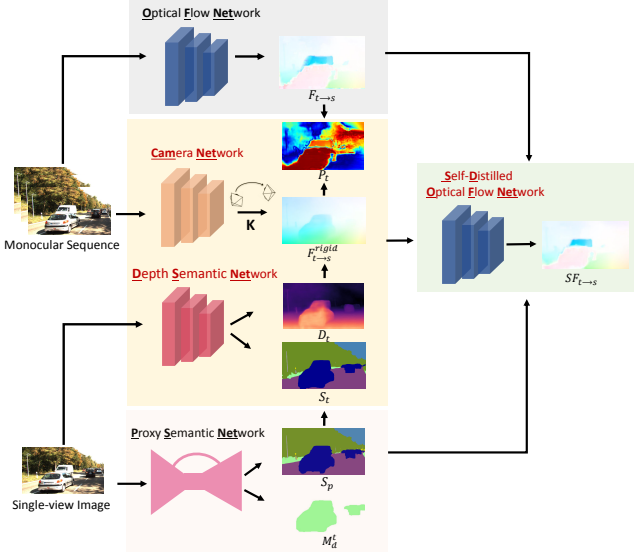


Figure 2. Overall framework for training Ω Net to predict depth, camera pose, camera intrinsics, semantic labels and optical flow. In red architectures composing Ω Net.

semantics of the scene. We can infer the depth of a scene by a single image mostly because of context and prior semantic knowledge. Prior works explored the possibility to learn both tasks with either full supervision [78, 23, 59, 45, 92, 44, 22] or supervision concerned with semantic labels only [90, 16]. Unlike previous works, we propose a compact architecture trained by self-supervision on monocular videos and exploiting proxy semantic labels.

Semantic segmentation and optical flow. Joint learning of semantic segmentation and optical flow estimation has been already explored [36]. Moreover, scene segmentation [72, 5] is required to disentangle potentially moving and static objects for focused optimizations. Differently, [66] leverages on optical flow to improve semantic predictions of moving objects. Peculiarly w.r.t. previous work, our proposal features a novel self-distillation training procedure guided by semantics to improve occlusion handling.

Scene understanding from stereo videos. Finally, we mention recent works approaching stereo depth estimation with optical flow [1] and semantic segmentation [42] for comprehensive scene understanding. In contrast, we are the first to rely on monocular videos to this aim.

3. Overall Learning Framework

Our goal is to develop a real-time comprehensive scene understanding framework capable of learning strictly related tasks from monocular videos. Purposely, we propose a multi-stage approach to learn first geometry and semantics, then elicit motion information, as depicted in Figure 2.

3.1. Geometry and Semantics

Self-supervised depth and pose estimation. We propose to solve a self-supervised single-image depth and pose estimation problem by exploiting geometrical constraints in a sequence of N images, in which one of the frames is used as the target view I_t and the other ones in turn as the source image I_s . Assuming a moving camera in a stationary scene, given a depth map D_t aligned with I_t , the camera intrinsic parameters K and the relative pose $T_{t \rightarrow s}$ between I_t and I_s , it is possible to sample pixels from I_s in order to synthesise a warped image \tilde{I}_t aligned with I_t . The mapping between corresponding homogeneous pixels coordinates $p_t \in I_t$ and $p_s \in I_s$ is given by:

$$p_s \sim K T_{t \rightarrow s} D_{p_t} K^{-1} p_t \quad (1)$$

Following [95], we use the sub-differentiable bilinear sampler mechanism proposed in [39] to obtain \tilde{I}_t . Thus, in order to learn depth, pose and camera intrinsics we train two separate CNNs to minimize the photometric reconstruction error between \tilde{I}_t and I_t , defined as:

$$\mathcal{L}_{ap}^D = \sum_p \psi(I_t(p), \tilde{I}_t(p)) \quad (2)$$

where ψ is a photometric error function between the two images. However, as pointed out in [29], such a formulation is prone to errors at occlusion/disocclusion regions or in static camera scenarios. To soften these issues, we follow the same principles as suggested in [29], where a minimum per-pixel reprojection loss is used to compute the photometric error, an automask method allows for filtering-out spurious gradients when the static camera assumption is violated, and an edge-aware smoothness loss term is used as in [28]. Moreover, we use the depth normalization strategy proposed in [77]. See supplementary material for further details.

We compute the rigid flow between I_t and I_s as the difference between the projected and original pixel coordinates in the target image:

$$F_{t \rightarrow s}^{rigid}(p_t) = p_s - p_t \quad (3)$$

Distilling semantic knowledge. The proposed distillation scheme is motivated by how time-consuming and cumbersome obtaining accurate pixel-wise semantic annotations is. Thus, we train our framework to estimate semantic segmentation masks S_t by means of supervision from cheap proxy labels S_p distilled by a semantic segmentation network, pre-trained on few annotated samples and capable to generalize well to diverse datasets. Availability of proxy semantic labels for the frames of a monocular video enables us to train a single network to predict jointly depth and semantic labels. Accordingly, the joint loss is obtained

by adding a standard cross-entropy term \mathcal{L}_{sem} to the previously defined self-supervised image reconstruction loss \mathcal{L}_{ap}^D . Moreover, similarly to [90], we deploy a cross-task loss term, \mathcal{L}_{edge}^D (see supplementary), aimed at favouring spatial coherence between depth edges and semantic boundaries. However, unlike [90], we do not exploit stereo pairs at training time.

3.2. Optical Flow and Motion Segmentation

Self-supervised optical flow. As the 3D structure of a scene includes stationary as well as non-stationary objects, to handle the latter we rely on a classical optical flow formulation. Formally, given two images I_t and I_s , the goal is to estimate the 2D motion vectors $F_{t \rightarrow s}(p_t)$ that map each pixel in I_t into its corresponding one in I_s . To learn such a mapping without supervision, previous approaches [58, 54, 88] employ an image reconstruction loss \mathcal{L}_{ap}^F that minimizes the photometric differences between I_t and the back-warped image \tilde{I}_t obtained by sampling pixels from I_s using the estimated 2D optical flow $F_{t \rightarrow s}(p_t)$. This approach performs well for non-occluded pixels but provides misleading information within occluded regions.

Pixel-wise motion probability. Non-stationary objects produce systematic errors when optimizing \mathcal{L}_{ap}^D due to the assumption that the camera is the only moving body in an otherwise stationary scene. However, such systematic errors can be exploited to identify non-stationary objects: at pixels belonging to such objects the rigid flow $F_{t \rightarrow s}^{rigid}$ and the optical flow $F_{t \rightarrow s}$ should exhibit different directions and/or norms. Therefore, a pixel-wise probability of belonging to an object independently moving between frames s and t , P_t , can be obtained by normalizing the differences between the two vectors. Formally, denoting with $\theta(p_t)$ the angle between the two vectors at location p_t , we define the per-pixel motion probabilities as:

$$P_t(p_t) = \max\left\{\frac{1 - \cos\theta(p_t)}{2}, 1 - \rho(p_t)\right\} \quad (4)$$

where $\cos\theta(p_t)$ can be computed as the normalized dot product between the vectors and evaluates the similarity in direction between them, while $\rho(p_t)$ is defined as

$$\rho(p_t) = \frac{\min\{\|F_{t \rightarrow s}(p_t)\|_2, \|F_{t \rightarrow s}^{rigid}(p_t)\|_2\}}{\max\{\|F_{t \rightarrow s}(p_t)\|_2, \|F_{t \rightarrow s}^{rigid}(p_t)\|_2\}}, \quad (5)$$

i.e. a normalized score of the similarity between the two norms. By taking the maximum of the two normalized differences, we can detect moving objects even when either the directions or the norms of the vectors are similar. A visualization of $P_t(p_t)$ is depicted in Fig. 3(d).

Semantic-aware Self-Distillation Paradigm. Finally, we combine semantic information, estimated optical flow, rigid flow and pixel-wise motion probabilities within a final

training stage to obtain a more robust self-distilled optical flow network. In other words, we train a new instance of the model to infer a self-distilled flow $SF_{t \rightarrow s}$ given the estimates $F_{t \rightarrow s}$ from a first self-supervised network and the aforementioned cues. As previously discussed and highlighted in Figure 3(c), standard self-supervised optical flow is prone to errors in occluded regions due to the lack of photometric information but can provide good estimates for the dynamic objects in the scene. On the contrary, the estimated rigid flow can properly handle occluded areas thanks to the minimum-reprojection mechanism [29]. Starting from these considerations, our key idea is to split the scene into stationary and potentially dynamics objects, and apply on them the proper supervision. Purposely, we can leverage several observations:

1. **Semantic priors.** Given a semantic map S_t for image I_t , we can binarize pixels into static M_t^s and potentially dynamic M_t^d , with $M_t^s \cap M_t^d = \emptyset$. For example, we expect that points labeled as *road* are static in the 3D world, while pixels belonging to the semantic class *car* may move. In M_t^d , we assign 1 for each potentially dynamic pixel, 0 otherwise, as shown in Figure 3(e).
2. **Camera Motion Boundary Mask.** Instead of using a backward-forward strategy [96] to detect boundaries occluded due to the ego-motion, we analytically compute a binary boundary mask M_t^b from depth and ego-motion estimates as proposed in [57]. We assign a 0 value for out-of-camera pixels, 1 otherwise as shown in Figure 3(f).
3. **Consistency Mask.** Because the inconsistencies between the rigid flow and $F_{t \rightarrow s}$ are not only due to dynamic objects but also to occluded/inconsistent areas, we can leverage Equation (4) to detect such critical regions. Indeed, we define the consistency mask as:

$$M_t^c = P_t < \xi, \xi \in [0, 1] \quad (6)$$

This mask assigns 1 where the condition is satisfied, 0 otherwise (*i.e.* inconsistent regions) as in Figure 3(g).

Finally, we compute the final mask M , in Figure 3(h), as:

$$M = \min\{\max\{M_t^d, M_t^c\}, M_t^b\} \quad (7)$$

As a consequence, M will effectively distinguish regions in the image for which we can not trust the supervision sourced by $F_{t \rightarrow s}$, *i.e.* inconsistent or occluded areas. On such regions, we can leverage our proposed self-distillation mechanism. Then, we define the final total loss for the self-distilled optical flow network as:

$$\mathcal{L} = \sum \alpha_r \phi(SF_{t \rightarrow s}, F_{t \rightarrow s}^{rigid}) \cdot (1 - M) + \alpha_d \phi(SF_{t \rightarrow s}, F_{t \rightarrow s}) \cdot M + \psi(I_t, \tilde{I}_t^{SF}) \cdot M \quad (8)$$

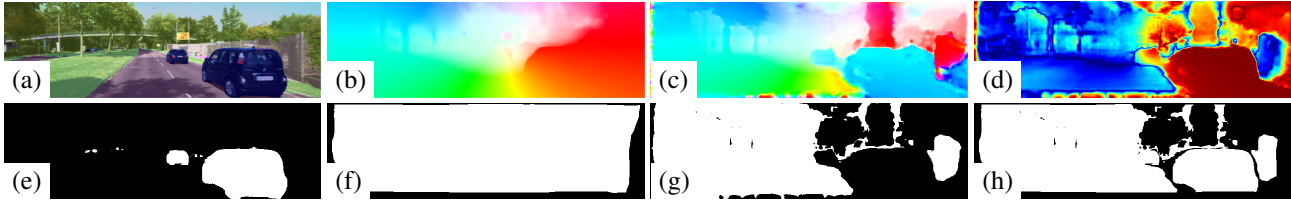


Figure 3. Overview of our semantic-aware and self-distilled optical flow estimation approach. We leverage semantic segmentation S_t (a) together with rigid flow $F_{t \rightarrow s}^{rigid}$ (b), teacher flow $F_{t \rightarrow s}$ (c) and motion probabilities P_t (d), the warmer the higher. From a) we obtain semantic priors M_t^d (e), combined with boundary mask M_t^b (f) and consistency mask M_t^c (g) derived from (d) as in Eq. 6, in order to obtain the final mask M (h) as in Eq. 7.

where ϕ is a distance function between two motion vectors, while α_r and α_d are two hyper-parameters.

3.3. Motion Segmentation

At test time, from pixel-wise probability P_t computed between $SF_{t \rightarrow s}$ and $F_{t \rightarrow s}^{rigid}$, semantic prior M_t^d and a threshold τ , we compute a motion segmentation mask by:

$$M_t^{mot} = M_t^d \cdot (P_t > \tau), \tau \in [0, 1] \quad (9)$$

Such mask allows us to detect moving objects in the scene independently of the camera motion. A qualitative example is depicted in Figure 1(f).

4. Architecture and Training Schedule

In this section we present the networks composing Ω Net (highlighted in red in Figure 2), and delineate their training protocol. We set $N = 3$, using 3-frames sequences. The source code is available at <https://github.com/CVLAB-Unibo/omeganet>.

4.1. Network architectures

We highlight the key traits of each network, referring the reader to the supplementary material for exhaustive details.

Depth and Semantic Network (DSNet). We build a single model, since shared reasoning about the two tasks is beneficial to both [90, 16]. To achieve real-time performance, DSNet is inspired to PydNet [62], with several key modifications due to the different goals. We extract a pyramid of features down to $\frac{1}{32}$ resolution, estimating a first depth map at the bottom. Then, it is upsampled and concatenated with higher level features in order to build a refined depth map. We repeat this procedure up to half resolution, where two estimators predict the final depth map D_t and semantic labels S_t . These are bi-linearly upsampled to full resolution. Each conv layer is followed by batch normalization and ReLU, but the prediction layers, using reflection padding. DSNet counts 1.93M parameters.

Camera Network (CamNet). This network estimates both camera intrinsics and poses between a target I_t and

some source views $I_s (1 \leq s \leq 3, s \neq t)$. CamNet differs from previous work by extracting features from I_t and I_s independently with shared encoders. We extract a pyramid of features down to $\frac{1}{16}$ resolution for each image and concatenate them to estimate the 3 Euler angles and the 3D translation for each I_s . As in [30], we also estimate the camera intrinsics. Akin to DSNet, we use batch normalization and ReLU after each layer but for prediction layers. CamNet requires 1.77M parameters for pose estimation and 1.02K for the camera intrinsics.

Optical Flow Network (OFNet). To pursue real-time performance, we deploy a 3-frame PWC-Net [75] network as in [54], which counts 4.79M parameters. Thanks to our novel training protocol leveraging on semantics and self-distillation, our OFNet can outperform other multi-task frameworks [3] built on the same optical flow architecture.

4.2. Training Protocol

Similarly to [88], we employ a two stage learning process to facilitate the network optimisation process. At first, we train DSNet and CamNet simultaneously, then we train OFNet by the self-distillation paradigm described in 3.2. For both stages, we use a batch size of 4 and resize input images to 640×192 for the KITTI dataset (and to 768×384 for pre-training on Cityscapes), optimizing the output of the networks at the highest resolution only. We also report additional experimental results for different input resolutions where specified. We use the Adam optimizer [46] with $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. As photometric loss ψ , we employ the same function defined in [28]. When training our networks, we apply losses using as I_s both the previous and the next image of our 3-frame sequence. Finally, we set both τ and ξ to be 0.5 in our experiments.

Depth, Pose, Intrinsics and Semantic Segmentation.

In order to train DSNet and CamNet we employ sequences of 3 consecutive frames and semantic proxy labels yielded by a state-of-the-art architecture [12] trained on Cityscapes with ground-truth labels. We trained DSNet and CamNet for 300K iterations, setting the initial learning rate to 10^{-4} , manually halved after 200K, 250K and 275K steps. We apply data augmentation to images as in [28]. Training takes

~ 20 hours on a Titan Xp GPU.

Optical Flow. We train OFNet by the procedure presented in 3.2. In particular, we perform 200K training steps with an initial learning rate of 10^{-4} , halved every 50K until convergence. Moreover, we apply strong data augmentation consisting in random horizontal and vertical flip, crops, random time order switch and, peculiarly, time stop, replacing all I_s with I_t to learn a zero motion vector. This configuration requires about 13 hours on a Titan Xp GPU with the standard 640×192 resolution. We use an L1 loss as ϕ . Once obtained a competitive network in non-occluded regions we train a more robust optical flow network, denoted as SD-OFNet, starting from pre-learned weights and the same structure of OFNet by distilling knowledge from OFNet and rigid flow computed by DSNet using the total mask M and 416×128 random crops applied to $F_{t \rightarrow s}$, $F_{t \rightarrow s}^{rigid}$, M and RGB images. We train SD-OFNet for 15K steps only with a learning rate of 2.5×10^{-5} halved after 5K, 7.5K, 10K and 12.5K steps, setting α_r to 0.025 and α_d to 0.2. At test-time, we rely on SD-OFNet only.

5. Experimental results

Using standard benchmark datasets, we present here the experimental validation on the main tasks tackled by Ω Net.

5.1. Datasets.

We conduct experiments on standard benchmarks such as KITTI and Cityscapes. We do not use feature extractors pre-trained on ImageNet or other datasets. For the sake of space, we report further studies in the supplementary material (*e.g.* results on pose estimation or generalization).

KITTI (K) [27] is a collection of 42,382 stereo sequences taken in urban environments from two video cameras and a LiDAR device mounted on the roof of a car. This dataset is widely used for benchmarking geometric understanding tasks such as depth, flow and pose estimation.

Cityscapes (CS) [19] is an outdoor dataset containing stereo pairs taken from a moving vehicle in various weather conditions. This dataset features higher resolution and higher quality images. While sharing similar settings, this dataset contains more dynamics scenes compared to KITTI. It consists of 22,973 stereo pairs with 2048×1024 resolution. 2,975 and 500 images come with fine semantic

5.2. Monocular Depth Estimation

In this section, we compare our results to other state-of-the-art proposals and assess the contribution of each component to the quality of our monocular depth predictions.

Comparison with state-of-the-art. We compare with state-of-the-art self-supervised networks trained on monocular videos according to the protocol described in [24]. We follow the same pre-processing procedure as [95] to remove

static images from the training split while using all the 697 images for testing. LiDAR points provided in [27] are re-projected on the left input image to obtain ground-truth labels for evaluation, up to 80 meters [25]. Since the predicted depth is defined up to a scale factor, we align the scale of our estimates by multiplying them by a scalar that matches the median of the ground-truth, as introduced in [95]. We adopt the standard performance metrics defined in [24]. Table 1 reports extensive comparison with respect to several monocular depth estimation methods. We outperform our main competitors such as [88, 96, 17, 3] that solve multi-task learning or other strategies that exploit additional information during the training/testing phase [9, 83]. Moreover, our best configuration, *i.e.* pre-training on CS and using 1024×320 resolution, achieves better results in 5 out of 7 metrics with respect to the single-task, state-of-the-art proposal [29] (and is the second best and very close to it on the remaining 2) which, however, leverages on a larger ImageNet pre-trained model based on ResNet-18. It is also interesting to note how our proposal without pre-training obtains the best performance in 6 out of 7 measures on 640×192 images (row 1 vs 15). These results validate our intuition about how the use of semantic information can guide geometric reasoning and make a compact network provide state-of-the-art performance even with respect to larger and highly specialized depth-from-mono methods.

Ablation study. Table 2 highlights how progressively adding the key innovations proposed in [30, 29, 77] contributes to strengthen Ω Net, already comparable to other methodologies even in its baseline configuration (first row). Interestingly, a large improvement is achieved by deploying joint depth and semantic learning (rows 5 vs 7), which forces the network to simultaneously reason about geometry and content within the same shared features. By replacing DSNet within Ω Net with a larger backbone [88] (rows 5 vs 6) we obtain worse performance, validating the design decisions behind our compact model. Finally, by pre-training on CS we achieve the best accuracy, which increases alongside with the input resolution (rows 8 to 10).

Depth Range Error Analysis. We dig into our depth evaluation to explain the effectiveness of Ω Net with respect to much larger networks. Table 3 compares, at different depth ranges, our model with more complex ones [29, 88]. This experiment shows how Ω Net superior performance comes from better estimation of large depths: Ω Net outperforms both competitors when we include distances larger than 8 m in the evaluation, while it turns out less effective in the close range.

5.3. Semantic Segmentation

In Table 4, we report the performance of Ω Net on semantic segmentation for the 19 evaluation classes of CS according to the metrics defined in [19, 4]. We compare Ω Net

Method	M	A	I	CS	Lower is better				Higher is better		
					Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Godard <i>et al.</i> [29]					0.132	1.044	5.142	0.210	0.845	0.948	0.977
Godard <i>et al.</i> [29] (1024 × 320)			✓		0.115	0.882	4.701	0.190	0.879	0.961	0.982
Zhou <i>et al.</i> [94]			✓		0.121	0.837	4.945	0.197	0.853	0.955	0.982
Mahjourian <i>et al.</i> [57]				✓	0.159	1.231	5.912	0.243	0.784	0.923	0.970
Wang <i>et al.</i> [77]				✓	0.151	1.257	5.583	0.228	0.810	0.936	0.974
Bian <i>et al.</i> [6]				✓	0.128	1.047	5.234	0.208	0.846	0.947	0.970
Yin <i>et al.</i> [88]	✓			✓	0.153	1.328	5.737	0.232	0.802	0.934	0.972
Zou <i>et al.</i> [96]	✓			✓	0.146	1.182	5.215	0.213	0.818	0.943	0.978
Chen <i>et al.</i> [17]	✓		✓		0.135	1.070	5.230	0.210	0.841	0.948	0.980
Luo <i>et al.</i> [56]	✓				0.141	1.029	5.350	0.216	0.816	0.941	0.976
Ranjan <i>et al.</i> [3]	✓				0.139	1.032	5.199	0.213	0.827	0.943	0.977
Xu <i>et al.</i> [83]		✓	✓		0.138	1.016	5.352	0.217	0.823	0.943	0.976
Casser <i>et al.</i> [9]		✓			0.141	1.026	5.290	0.215	0.816	0.945	0.979
Gordon <i>et al.</i> [30]	✓	✓			0.128	0.959	5.230	-	-	-	-
Ω Net(640 × 192)	✓	✓			0.126	0.835	4.937	0.199	0.844	0.953	0.982
Ω Net(1024 × 320)	✓	✓			0.125	0.805	4.795	0.195	0.849	0.955	0.983
Ω Net(640 × 192)	✓	✓		✓	0.120	0.792	4.750	0.191	0.856	0.958	0.984
Ω Net(1024 × 320)	✓	✓		✓	0.118	0.748	4.608	0.186	0.865	0.961	0.985

Table 1. Depth evaluation on the Eigen split [24] of KITTI [26]. We indicate additional features of each method. M: multi-task learning, A: additional information (e.g. object knowledge, semantic information), I: feature extractors pre-trained on ImageNet [20], CS: network pre-trained on Cityscapes [19].

Resolution	Learned Intr. [30]	Norm. [77]	Min. Repr. [29]	Automask [29]	Sem. [12]	Pre-train	Lower is better				Higher is better		
							Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
640 × 192	-	-	-	-	-	-	0.139	1.056	5.288	0.215	0.826	0.942	0.976
640 × 192	✓	-	-	-	-	-	0.138	1.014	5.213	0.213	0.829	0.943	0.977
640 × 192	✓	✓	-	-	-	-	0.136	1.008	5.204	0.212	0.832	0.944	0.976
640 × 192	✓	✓	✓	-	-	-	0.132	0.960	5.104	0.206	0.840	0.949	0.979
640 × 192	✓	✓	✓	✓	-	-	0.130	0.909	5.022	0.207	0.842	0.948	0.979
640 × 192 †	✓	✓	✓	✓	-	-	0.134	1.074	5.451	0.213	0.834	0.946	0.977
640 × 192	✓	✓	✓	✓	✓	-	0.126	0.835	4.937	0.199	0.844	0.953	0.980
416 × 128	✓	✓	✓	✓	✓	✓	0.126	0.862	4.963	0.199	0.846	0.952	0.981
640 × 192	✓	✓	✓	✓	✓	✓	0.120	0.792	4.750	0.191	0.856	0.958	0.984
1024 × 320	✓	✓	✓	✓	✓	✓	0.118	0.748	4.608	0.186	0.865	0.961	0.985

Table 2. Ablation study of our depth network on the Eigen split [24] of KITTI. †: our network is replaced by a ResNet50 backbone [88].

Method	Cap (m)	Abs Rel	Sq Rel	RMSE	RMSE log
Godard <i>et al.</i> [29]	0-8	0.059	0.062	0.503	0.082
Ω Net†	0-8	0.060	0.063	0.502	0.082
Ω Net	0-8	0.062	0.065	0.517	0.085
Godard <i>et al.</i> [29]	0-50	0.125	0.788	3.946	0.198
Ω Net†	0-50	0.127	0.762	4.020	0.199
Ω Net	0-50	0.124	0.702	3.836	0.195
Godard <i>et al.</i> [29]	0-80	0.132	1.044	5.142	0.210
Ω Net†	0-80	0.134	1.074	5.451	0.213
Ω Net	0-80	0.126	0.835	4.937	0.199

Table 3. Depth errors by varying the range. †: our network is replaced by a ResNet50 backbone [88].

Method	Train	Test	mIoU Class	mIoU Cat.	Pix.Acc.
DABNet [50]	CS(S)	CS	69.62	87.56	94.62
FCHardNet [11]	CS(S)	CS	76.37	89.22	95.35
Ω Net	CS(P)	CS	54.80	82.92	92.50
DABNet [50]	CS(S)	K	35.40	61.49	80.50
FCHardNet [11]	CS(S)	K	44.74	68.20	72.07
Ω Net	CS(P)	K	43.80	74.31	88.31
Ω Net	CS(P) + K(P)	K	46.68	75.84	88.12

Table 4. Semantic segmentation on Cityscapes (CS) and KITTI (K). S: training on ground-truth, P: training on proxy labels.

against state-of-the-art networks for real-time semantic segmentation [11, 50] when training on CS and testing either on the validation set of CS (rows 1-3) or the 200 semantically

annotated images of K (rows 4-6). Even though our network is not as effective as the considered methods when training and testing on the same dataset, it shows greater generalization capabilities to unseen domains: it significantly outperforms other methods when testing on K for mIoU_{category} and pixel accuracy, and provides similar results to [11] for mIoU_{class}. We relate this ability to our training protocol based on proxy labels (P) instead of ground truths (S). We validate this hypothesis with thorough ablation studies reported in the supplementary material. Moreover, as we have already effectively distilled the knowledge from DPC [12] during pre-training on CS, there is only a slight benefit in training on both CS and K (with proxy labels only) and testing on K (row 7). Finally, although achieving 46.68 mIoU on fine segmentation, we obtain 89.64 mIoU for the task of segmenting static from potentially dynamic classes, an important result to obtain accurate motion masks.

5.4. Optical Flow

In Table 5, we compare the performance of our optical flow network with competing methods using the KITTI 2015 stereo/flow training set [26] as testing set, which contains 200 ground-truth optical flow measurements for eval-

Method	Dataset	train			test
		Noc	All	F1	F1
Meister <i>et al.</i> [58] - C	SYN + K	-	8.80	28.94%	29.46%
Meister <i>et al.</i> [58] - CSS	SYN + K	-	8.10	23.27%	23.30%
Zou <i>et al.</i> [96]	SYN + K	-	8.98	26.0%	25.70%
Ranjan <i>et al.</i> [3]	SYN + K	-	5.66	20.93%	25.27%
Wang <i>et al.</i> [79] **	K	-	5.58	-	18.00%
Yin <i>et al.</i> [88]	K	8.05	10.81	-	-
Chen <i>et al.</i> [17] †	K	5.40	8.95	-	-
Chen <i>et al.</i> [17] (online) †	K	4.86	8.35	-	-
Ranjan <i>et al.</i> [3]	K	-	6.21	26.41%	-
Luo <i>et al.</i> [56]	K	-	5.84	-	21.56%
Luo <i>et al.</i> [56] *	K	-	5.43	-	20.61%
Ω Net (Ego-motion)	K	11.72	13.50	51.22%	-
OFNet	K	3.48	11.61	25.78%	-
SD-OFNet	K	3.29	5.39	20.0%	19.47%

Table 5. Optical flow evaluation on the KITTI 2015 dataset. †: pre-trained on ImageNet, SYN: pre-trained on SYNTHIA [71], *: trained on stereo pairs, **: using stereo at testing time.

uation. We exploit all the raw K images for training, but we exclude the images used at testing time as done in [96], to be consistent with experimental results of previous self-supervised optical flow strategies [88, 96, 17, 3]. From the table, we can observe how our self-distillation strategy allows SD-OFNet to outperform by a large margin competitors trained on K only (rows 5-11), and it even performs better than models pre-initialized by training on synthetic datasets [71]. Moreover, we submitted our flow predictions to the online KITTI flow benchmark after retraining the network including images from the whole official training set. In this configuration, we can observe how our model achieves state-of-the-art F1 performances with respect to other monocular multi-task architectures.

5.5. Motion Segmentation

In Table 6 we report experimental results for the motion segmentation task on the KITTI 2015 dataset, which provides 200 images manually annotated with motion labels for the evaluation. We compare our methodology with respect to other state-of-the-art strategies that performs multi-task learning and motion segmentation [3, 56, 79] using the metrics and evaluation protocol proposed in [56]. It can be noticed how our segmentation strategy outperforms all the other existing methodologies by a large margin. This demonstrates the effectiveness of our proposal to jointly combine semantic reasoning and motion probability to obtain much better results. We also report, as upper bound, the accuracy enabled by injecting semantic proxies [12] in place of Ω Net semantic predictions to highlight the low margin between the two.

5.6. Runtime analysis

Finally, we measure the runtime of Ω Net on different hardware devices, *i.e.* a Titan Xp GPU, an embedded NVIDIA Jetson TX2 board and an Intel i7-7700K@4.2 GHz CPU. Timings averaged over 200 frames at 640×192

Method	Pixel Acc.	Mean Acc.	Mean IoU	f.w. IoU
Yang <i>et al.</i> [86] *	0.89	0.75	0.52	0.87
Luo <i>et al.</i> [56]	0.88	0.63	0.50	0.86
Luo <i>et al.</i> [56] *	0.91	0.76	0.53	0.87
Wang <i>et al.</i> [79] (Full) **	0.90	0.82	0.56	0.88
Ranjan <i>et al.</i> [3]	0.87	0.79	0.53	0.85
ΩNet	0.98	0.86	0.75	0.97
Ω Net (Proxy [12])	0.98	0.87	0.77	0.97

Table 6. Motion segmentation evaluation on the KITTI 2015 dataset. *: trained on stereo pairs, **: using stereo at testing time.

Device	Watt	D	DS	OF	Cam	Ω
Jetson TX2	15	12.5	10.3	6.5	49.2	4.5
i7-7700K	91	5.0	4.2	4.9	31.4	2.4
Titan XP	250	170.2	134.1	94.1	446.7	57.4

Table 7. Runtime analysis on different devices. We report the power consumption in Watt and the FPS. D: Depth, S: Semantic, OF: Optical Flow, Cam: camera pose, Ω : Overall architecture.

resolution. Moreover, as each component of Ω Net may be used on its own, we report the runtime for each independent task. As summarized in Table 7, our network runs in real-time on the Titan Xp GPU and at about 2.5 FPS on a standard CPU. It also fits the low-power NVIDIA Jetson TX2, achieving 4.5 FPS to compute all the outputs. Additional experiments are available in the supplementary material.

6. Conclusions

In this paper, we have proposed the first real-time network for comprehensive scene understanding from monocular videos. Our framework reasons jointly about geometry, motion and semantics in order to estimate accurately depth, optical flow, semantic segmentation and motion masks at about 60 FPS on high-end GPU and 5FPS on embedded systems. To address the above multi-task problem we have proposed a novel learning procedure based on distillation of proxy semantic labels and semantic-aware self-distillation of optical-flow information. Thanks to this original paradigm, we have demonstrated state-of-the-art performance on standard benchmark datasets for depth and optical flow estimation as well as for motion segmentation.

As for future research, we find it intriguing to investigate on whether and how would it be possible to self-adapt Ω Net on-line. Although some very recent works have explored this topic for depth-from-mono [9] and optical flow [17], the key issue with our framework would be to conceive a strategy to deal with semantics.

Acknowledgement. We gratefully acknowledge the support of NVIDIA Corporation with the donation of the Titan X Pascal GPU used for this research.

References

- [1] Filippo Aleotti, Matteo Poggi, Fabio Tosi, and Stefano Mattochia. Learning end-to-end scene flow by distilling single

- tasks knowledge. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020. 3
- [2] Lorenzo Andraghetti, Panteleimon Myriokefalitakis, Pier Luigi Dovesi, Belen Luque, Matteo Poggi, Alessandro Pieropan, and Stefano Mattoccia. Enhancing self-supervised monocular depth estimation with traditional visual odometry. In *7th International Conference on 3D Vision (3DV)*, 2019. 2
- [3] Ranjan Anurag, Varun Jampani, Kihwan Kim, Deqing Sun, Jonas Wulff, and Michael J. Black. Competitive collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2, 5, 6, 7, 8
- [4] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017. 2, 6
- [5] Min Bai, Wenjie Luo, Kaustav Kundu, and Raquel Urtasun. Exploiting semantic information and deep matching for optical flow. In *European Conference on Computer Vision*, pages 154–170. Springer, 2016. 3
- [6] Jia-Wang Bian, Zhichao Li, Naiyan Wang, Huangying Zhan, Chunhua Shen, Ming-Ming Cheng, and Ian Reid. Unsupervised scale-consistent depth and ego-motion learning from monocular video. In *Thirty-third Conference on Neural Information Processing Systems (NeurIPS)*, 2019. 2, 7
- [7] Michael J Black and Paul Anandan. The robust estimation of multiple motions: Parametric and piecewise-smooth flow fields. *Computer vision and image understanding*, 63(1):75–104, 1996. 2
- [8] Thomas Brox, Andrés Bruhn, Nils Papenberg, and Joachim Weickert. High accuracy optical flow estimation based on a theory for warping. In *European conference on computer vision*, pages 25–36. Springer, 2004. 2
- [9] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In *Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)*, 2019. 2, 6, 7, 8
- [10] Ya-Liang Chang, Zhe Yu Liu, and Winston Hsu. Vornet: Spatio-temporally consistent video inpainting for object removal. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 2
- [11] Ping Chao, Chao-Yang Kao, Yu-Shan Ruan, Chien-Hsiang Huang, and Youn-Long Lin. Hardnet: A low memory traffic network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3552–3561, 2019. 2, 7
- [12] Liang-Chieh Chen, Maxwell Collins, Yukun Zhu, George Papandreou, Barret Zoph, Florian Schroff, Hartwig Adam, and Jon Shlens. Searching for efficient multi-scale architectures for dense image prediction. In *Advances in Neural Information Processing Systems*, pages 8699–8710, 2018. 2, 5, 7, 8
- [13] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 2
- [14] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. 2
- [15] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018. 2
- [16] Po-Yi Chen, Alexander H. Liu, Yen-Cheng Liu, and Yu-Chiang Frank Wang. Towards scene understanding: Unsupervised monocular depth estimation with semantic-aware representation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 3, 5
- [17] Yuhua Chen, Cordelia Schmid, and Cristian Sminchisescu. Self-supervised learning with geometric constraints in monocular video: Connecting flow, depth, and camera. In *ICCV*, 2019. 2, 6, 7, 8
- [18] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017. 2
- [19] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 6, 7
- [20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 7
- [21] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, 2015. 2
- [22] Pier Luigi Dovesi, Matteo Poggi, Lorenzo Andraghetti, Miquel Martí, Hedvig Kjellström, Alessandro Pieropan, and Stefano Mattoccia. Real-time semantic stereo matching. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2020. 3
- [23] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE international conference on computer vision*, pages 2650–2658, 2015. 3
- [24] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374, 2014. 2, 6, 7
- [25] Ravi Garg, Vijay Kumar BG, Gustavo Carneiro, and Ian Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *European Conference on Computer Vision*, pages 740–756. Springer, 2016. 2, 6

- [26] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 7
- [27] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *CVPR*, 2012. 6
- [28] Clément Godard, Oisín Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, 2017. 2, 3, 5
- [29] Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *ICCV*, 2019. 2, 3, 4, 6, 7
- [30] Ariel Gordon, Hanhan Li, Rico Jonschkowski, and Anelia Angelova. Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. In *ICCV*, 2019. 2, 5, 6, 7
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 37(9):1904–1916, 2015. 2
- [32] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2
- [33] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):185–203, 1981. 2
- [34] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. 2
- [35] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. Lite-flownet: A lightweight convolutional neural network for optical flow estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [36] Junhwa Hur and Stefan Roth. Joint optical flow and temporally consistent semantic segmentation. In *European Conference on Computer Vision*, pages 163–177. Springer, 2016. 3
- [37] Junhwa Hur and Stefan Roth. Iterative residual refinement for joint optical flow and occlusion estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5754–5763, 2019. 2
- [38] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [39] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in neural information processing systems*, pages 2017–2025, 2015. 3
- [40] Joel Janai, Fatma G’uney, Anurag Ranjan, Michael J. Black, and Andreas Geiger. Unsupervised learning of multi-frame optical flow with occlusions. In *European Conference on Computer Vision (ECCV)*, volume Lecture Notes in Computer Science, vol 11220, pages 713–731. Springer, Cham, Sept. 2018. 2
- [41] J Yu Jason, Adam W Harley, and Konstantinos G Derpanis. Back to basics: Unsupervised learning of optical flow via brightness constancy and motion smoothness. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. 2
- [42] Huaizu Jiang, Deqing Sun, Varun Jampani, Zhaoyang Lv, Erik Learned-Miller, and Jan Kautz. Sense: A shared encoder network for scene-flow estimation. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 3
- [43] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super slomo: High quality estimation of multiple intermediate frames for video interpolation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [44] Jianbo Jiao, Ying Cao, Yibing Song, and Rynson Lau. Look deeper into depth: Monocular depth estimation with semantic booster and attention-driven loss. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 53–69, 2018. 3
- [45] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weight losses for scene geometry and semantics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7482–7491, 2018. 3
- [46] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [47] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 2
- [48] Yevhen Kuznetsov, Jorg Stuckler, and Bastian Leibe. Semi-supervised deep learning for monocular depth map prediction. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 2
- [49] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. In *3DV*, 2016. 2
- [50] Gen Li and Joongkyu Kim. Dabnet: Depth-wise asymmetric bottleneck for real-time semantic segmentation. In *British Machine Vision Conference*, 2019. 2, 7
- [51] Hanchao Li, Pengfei Xiong, Haoqiang Fan, and Jian Sun. Dfanet: Deep feature aggregation for real-time semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9522–9531, 2019. 2
- [52] Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, Wei Hua, Alan Yuille, and Li Fei-Fei. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation. In *CVPR*, 2019. 2
- [53] Pengpeng Liu, Irwin King, Michael R. Lyu, and Jia Xu. Ddfow: Learning optical flow with unlabeled data distillation. In *AAAI*, 2019. 2
- [54] Pengpeng Liu, Michael R. Lyu, Irwin King, and Jia Xu. Self-low: Self-supervised learning of optical flow. In *CVPR*, 2019. 2, 4, 5

- [55] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 2
- [56] Chenxu Luo, Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, Ram Nevatia, and Alan Yuille. Every pixel counts++: Joint learning of geometry and motion with 3d holistic understanding. *PAMI*, 2019. 2, 7, 8
- [57] Reza Mahjourian, Martin Wicke, and Anelia 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)*, 2018. 2, 4, 7
- [58] Simon Meister, Junhwa Hur, and Stefan Roth. UnFlow: Unsupervised learning of optical flow with a bidirectional census loss. In *AAAI*, New Orleans, Louisiana, Feb. 2018. 1, 2, 4, 8
- [59] Arsalan Mousavian, Hamed Pirsiavash, and Jana Košecák. Joint semantic segmentation and depth estimation with deep convolutional networks. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 611–619. IEEE, 2016. 3
- [60] Adam Paszke, Abhishek Chaurasia, Sangpil Kim, and Eugenio Culurciello. Enet: A deep neural network architecture for real-time semantic segmentation. *arXiv preprint arXiv:1606.02147*, 2016. 2
- [61] Valentino Peluso, Antonio Cipolletta, Andrea Calimera, Matteo Poggi, Fabio Tosi, and Stefano Mattoccia. Enabling energy-efficient unsupervised monocular depth estimation on ARMv7-based platforms. In *Design, Automation & Test in Europe Conference & Exhibition, DATE 2019, Florence, Italy, March 25-29, 2019*, pages 1703–1708, 2019. 2
- [62] Matteo Poggi, Filippo Aleotti, Fabio Tosi, and Stefano Mattoccia. Towards real-time unsupervised monocular depth estimation on CPU. In *IEEE/RSJ Conference on Intelligent Robots and Systems (IROS)*, 2018. 2, 5
- [63] Matteo Poggi, Filippo Aleotti, Fabio Tosi, and Stefano Mattoccia. On the uncertainty of self-supervised monocular depth estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [64] Matteo Poggi, Fabio Tosi, and Stefano Mattoccia. Learning monocular depth estimation with unsupervised trinocular assumptions. In *6th International Conference on 3D Vision (3DV)*, 2018. 2
- [65] Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [66] Hazem Rashed, Senthil Yogamani, Ahmad El-Sallab, Pavel Krizek, and Mohamed El-Helw. Optical flow augmented semantic segmentation networks for automated driving. *arXiv preprint arXiv:1901.07355*, 2019. 3
- [67] Zhile Ren, Orazio Gallo, Deqing Sun, Ming-Hsuan Yang, Erik Sudderth, and Jan Kautz. A fusion approach for multi-frame optical flow estimation. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 2077–2086. IEEE, 2019. 2
- [68] Zhe Ren, Junchi Yan, Bingbing Ni, Bin Liu, Xiaokang Bin, and Hongyuan Zha. Unsupervised deep learning for optical flow estimation. In *Thirty-First AAAI Conference on Artificial Intelligence (AAAI-17)*, 2017. 2
- [69] J. Revaud, P Weinzaepfel, Z. Harchaoui, and C. Schmid. Epicflow: Edge-preserving interpolation of correspondences for optical flow. In *CVPR*, 2019. 2
- [70] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 2
- [71] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio Lopez. The SYNTHIA Dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *CVPR*, 2016. 8
- [72] Laura Sevilla-Lara, Deqing Sun, Varun Jampani, and Michael J Black. Optical flow with semantic segmentation and localized layers. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3889–3898, 2016. 3
- [73] Wei-Shi Zheng Shuosun Guan, Haoxin Li. Unsupervised learning for optical flow estimation using pyramid convolution lstm. In *Proceedings of IEEE International Conference on Multimedia and Expo(ICME)*, 2019. 2
- [74] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 2
- [75] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 5
- [76] Fabio Tosi, Filippo Aleotti, Matteo Poggi, and Stefano Mattoccia. Learning monocular depth estimation infusing traditional stereo knowledge. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [77] Chaoyang Wang, Jose Miguel Buenaposada, Rui Zhu, and Simon Lucey. Learning depth from monocular videos using direct methods. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 3, 6, 7
- [78] Peng Wang, Xiaohui Shen, Zhe Lin, Scott Cohen, Brian Price, and Alan L Yuille. Towards unified depth and semantic prediction from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2800–2809, 2015. 3
- [79] Yang Wang, Peng Wang, Zhenheng Yang, Chenxu Luo, Yi Yang, and Wei Xu. Unos: Unified unsupervised optical-flow and stereo-depth estimation by watching videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8071–8081, 2019. 8
- [80] Jamie Watson, Michael Firman, Gabriel J Brostow, and Daniyar Turmukhambetov. Self-supervised monocular depth hints. In *ICCV*, 2019. 2
- [81] Wofk, Diana and Ma, Fangchang and Yang, Tien-Ju and Karaman, Sertac and Sze, Vivienne. FastDepth: Fast Monocular Depth Estimation on Embedded Systems. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2019. 2

- [82] Yu Xiang, Alexandre Alahi, and Silvio Savarese. Learning to track: Online multi-object tracking by decision making. In *Proceedings of the IEEE international conference on computer vision*, pages 4705–4713, 2015. 2
- [83] Haoifei Xu, Jianmin Zheng, Jianfei Cai, and Juyong Zhang. Region deformer networks for unsupervised depth estimation from unconstrained monocular videos. In *IJCAI*, 2019. 2, 6, 7
- [84] Nan Yang, Rui Wang, Jörg Stückler, and Daniel 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
- [85] Wang Yang, Yi Yang, Zhenheng Yang, Liang Zhao, and Wei Xu. Occlusion aware unsupervised learning of optical flow. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [86] Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, and Ram Nevatia. Every pixel counts: Unsupervised geometry learning with holistic 3d motion understanding. In *The European Conference on Computer Vision (ECCV) Workshops*, September 2018. 2, 8
- [87] Zhenheng Yang, Peng Wang, Wang Yang, Wei Xu, and Nevatia Ram. Lego: Learning edge with geometry all at once by watching videos. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [88] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 2, 4, 5, 6, 7, 8
- [89] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 325–341, 2018. 2
- [90] Pierluigi Zama Ramirez, Matteo Poggi, Fabio Tosi, Stefano Mattoccia, and Luigi Di Stefano. Geometry meets semantic for semi-supervised monocular depth estimation. In *14th Asian Conference on Computer Vision (ACCV)*, 2018. 1, 3, 4, 5
- [91] Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian 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)*, 2018. 2
- [92] Zhenyu Zhang, Zhen Cui, Chunyan Xu, Zequn Jie, Xiang Li, and Jian Yang. Joint task-recursive learning for semantic segmentation and depth estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 235–251, 2018. 3
- [93] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017. 2
- [94] Junsheng Zhou, Yuwang Wang, Naiyan Wang, and Wenjun Zeng. Unsupervised high-resolution depth learning from videos with dual networks. In *Inter. Conf. on Computer Vision*. IEEE, IEEE, October 2019. 2, 7
- [95] Tinghui Zhou, Matthew Brown, Noah Snavely, and David 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, 3, 6
- [96] Yuliang Zou, Zelun Luo, and Jia-Bin Huang. Df-net: Unsupervised joint learning of depth and flow using cross-task consistency. In *European Conference on Computer Vision (ECCV)*, 2018. 1, 2, 4, 6, 7, 8

Distilled Semantics for Comprehensive Scene Understanding from Videos – Supplementary material

Fabio Tosi* Filippo Aleotti* Pierluigi Zama Ramirez*
Matteo Poggi Samuele Salti Luigi Di Stefano Stefano Mattoccia
Department of Computer Science and Engineering (DISI)
University of Bologna, Italy

*{fabio.tosi5, filippo.aleotti2, pierluigi.zama}@unibo.it

Supplementary material

This document provides additional material concerning CVPR 2020 paper, “Distilled Semantics for Comprehensive Scene Understanding from Videos”. In particular, we report here a more detailed description of our Ω Net architecture and the losses used to train it, alongside with more insights related to performance in the addressed tasks (depth, pose, optical flow, semantic and motion segmentation) and runtime. Moreover, we include additional qualitative results on KITTI (K) and CityScapes (CS), as well as on an arbitrary YouTube video for which the camera parameters are not known in advance, thus showing how Ω Net can provide comprehensive scene understanding in the wild.

1. Network Architecture

In this section, we provide a more detailed description of our Ω Net architecture.

Table 1 reports a detailed specification of the layers building up the DSNet and CamNet modules. For each layer, we report kernel size (K), stride (S) and number of input/output channels. As for OFNet and the proxy semantic network, a thorough description can be found in [17] and [6] respectively.

2. Losses

To train the DSNet module, we rely on a multi-task loss function based mainly on two terms. In particular, a depth term is in charge of minimize the discrepancy between the target image I_t and an image I_s , warped as \tilde{I}_t^s , from a monocular sequence while a semantic term is used to learn semantic labels from proxy label distilled by a pre-trained network.

Depth term. According to the self-supervised training paradigm proposed in [13], we adopt a photometric loss function consisting in a weighted combination between the Structural Dissimilarity Measure (DSSIM) and the standard \mathcal{L}_1 loss. In addition, a per-pixel minimum strategy [14] is used to solve occlusion/disocclusion by simply picking the minimum error between each pair I_t and I_s instead of averaging them. Thus, the photometric loss function is defined as:

$$\mathcal{L}_{ap}^D = \sum_p \min_s (\alpha \mathcal{L}_{DSSIM}(p) + (1 - \alpha) \|I_t(p) - \tilde{I}_t^s(p)\|_1) \quad (1)$$

where p denotes pixel coordinates, \tilde{I}_t^s a source image I_s warped according to estimated depth and pose and the $DSSIM$ loss function is computed as:

$$\mathcal{L}_{DSSIM}(p) = \frac{1 - SSIM(I_t(p), \tilde{I}_t^s(p))}{2} \quad (2)$$

In our experiments, we set $\alpha = 0.85$.

*Joint first authorship.

Layer	K	S	In/Out	Input
Deep feature extractor (DSE)				
conv1a	3	2	3/16	input
conv1b	3	1	16/16	conv1a
conv2a	3	2	16/32	conv1b
conv2b	3	1	32/32	conv2a
conv3a	3	2	32/64	conv2b
conv3b	3	1	64/64	conv3a
conv4a	3	2	64/128	conv3b
conv4b	3	1	128/128	conv4a
conv5a	3	2	128/256	conv4b
conv5b	3	1	256/256	conv5a
Estimator (E)				
conv1	3	1	i_channels/64	input
conv2	3	1	64/48	conv1
conv3	3	1	48/32	conv2
conv4	3	1	32/16	conv3
Context (C)				
disp1	3	1	i_channels/64	input
disp2	3	1	64/32	disp1
disp3	3	1	32/16	disp2
disp	3	1	16/1	disp2
Disparity and Semantic Tower				
conv5	3	1	i_channels/16	E(conv5b)
disp5	3	1	i_channels//1	C(conv5)
conv4	3	1	i_channels/16	E(conv4b, disp5 ↑)
disp4	3	1	i_channels//1	C(conv4, conv5 ↑) + disp5 ↑
conv3	3	1	i_channels/16	E(conv3b, disp4 ↑)
disp3	3	1	i_channels//1	C(conv3, conv4 ↑) + disp4 ↑
conv2	3	1	i_channels/16	E(conv2b, disp3 ↑)
disp2	3	1	i_channels//1	C(conv2, conv3 ↑) + disp3 ↑
conv1	3	1	i_channels/16	E(conv1b, disp2 ↑)
disp1	3	1	i_channels//1	C(conv1, conv2 ↑) + disp2 ↑
sem	3	1	i_channels//1	C(conv1, conv2 ↑) + disp2 ↑

Layer	K	S	In/Out	Input
Deep feature extractor (DFE)				
conv1a	3	2	3/16	input
conv1b	3	1	16/16	conv1a
conv2a	3	2	16/32	conv1b
conv2b	3	1	32/32	conv2a
conv3a	3	2	32/64	conv2b
conv3b	3	1	64/64	conv3a
conv4a	3	2	64/128	conv3b
conv4b	3	1	128/128	conv4a
Pose Estimator				
conv1a	3	1	i_channels/128	DFE_t, DFE_s
conv1b	3	2	128/128	conv1a
conv2a	3	1	128/256	conv1b
conv2b	3	2	256/256	conv2a
pose	1	1	256/6*N	conv2b
Intrinsic Estimator				
focals	1	1	i_channels/2	conv2b
offsets	1	1	i_channels/2	conv2b

Table 1. Detailed structure of the DSNet (left) and CamNet (right) modules in Ω Net. The symbol “;” means concatenation, while \uparrow indicates upsampling.

A smoothness term is also used to penalize large disparity differences between adjacent pixels when the former do not co-occur with strong RGB gradients:

$$\mathcal{L}_{smooth} = \sum_p |\nabla D_t(p)| \cdot \left(e^{-|\nabla I_t(p)|} \right)^T \quad (3)$$

Finally, we mask-out pixels whose appearance do not change between consecutive frames, which includes scenes with no relative motion. This has the effect of letting the network ignore pixels which move at the same velocity as the camera, and even to ignore whole frames when the camera stop moving. According to [14], this is accomplished by removing those pixels which have an unwarped photometric loss smaller than the corresponding warped photometric loss, *i.e.*

$$\mu = \min \mathcal{L}_{ap}^D(I_t, I_s) > \mathcal{L}_{ap}^D(I_t, I_t^s) \quad (4)$$

Semantic term. The standard cross-entropy loss between the predicted and proxy pixel-wise semantic labels is used as semantic term:

$$\mathcal{L}_{sem} = -(S_t \log(S_P) + (1 - S_t) \log(1 - S_P)) \quad (5)$$

where S_t is the semantics predicted by DSNet and S_P the ground-truth proxy label. Moreover, as proposed in [24] we employ a cross-task loss to tighten the link between the learning tasks dealing with depth and semantics:

$$\mathcal{L}_{cdd} = \sum_p \text{sgn}(|\nabla S_t(p)|) \cdot \left(e^{-|\nabla D_t(p)|} \right)^T \quad (6)$$

Hence, the total loss used to train DSNet is a weighted combination of the above losses:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{ap}^D + \lambda_2 \mathcal{L}_{smooth} + \lambda_3 \mathcal{L}_{sem} + \lambda_4 \mathcal{L}_{cdd} \quad (7)$$

where $\lambda_1, \lambda_2, \lambda_3$ and λ_4 are hyper-parameters. In our experiments, we set $\lambda_1 = 1, \lambda_2 = 0.1, \lambda_3 = 1$ and $\lambda_4 = 0.1$.

As described in the paper, for the Optical Flow we rely on a peculiar training schedule based on two components in Ω Net, which are referred as OFNet and SD-OFNet.

Optical Flow term. We train a the first instance of the optical flow network (OFNet) using the same photometric loss as for DSNet:

$$\mathcal{L}_{ap}^{OF} = \sum_p \alpha \mathcal{L}_{DSSIM} + (1 - \alpha) \|I_t - \tilde{I}_t^s\|_1 \quad (8)$$

In this case, however, \tilde{I}_t^s is warped according to estimated flow. Akin to DSNet, we set $\alpha = 0.85$.

Self-Distilled Optical Flow term. The self-distilled optical flow network (SD-OFNet), instead, is trained in a quite different manner. In fact, given the optical flow $F_{t \rightarrow s}$ predicted by OFNet, the rigid flow $F_{t \rightarrow s}^{rigid}$ and the mask M , we leverage on the optical flow in the regions where $F_{t \rightarrow s}$ and $F_{t \rightarrow s}^{rigid}$ are similar as well as on moving objects, while we rely on the rigid flow for the remaining areas (*e.g.*, occlusions due to camera motion). We can distinguish the former regions from the latter ones looking at M . Moreover, we also apply a photometric term ϕ on the predicted optical flow $SF_{t \rightarrow s}$. The final loss \mathcal{L} to train SD-OFNet is given by:

$$\mathcal{L} = \sum \alpha_r \phi(SF_{t \rightarrow s}, F_{t \rightarrow s}^{rigid}) \cdot (1 - M) + \alpha_d \phi(SF_{t \rightarrow s}, F_{t \rightarrow s}) \cdot M + \psi(I_t, \tilde{I}_t^{SF}) \cdot M \quad (9)$$

During training, $F_{t \rightarrow s}, F_{t \rightarrow s}^{rigid}, M$ and the input images of SD-OFNet are randomly cropped to 416×128 before computing \mathcal{L} : in doing so, the errors at occluded areas in $F_{t \rightarrow s}$ due to camera motions, clearly visible in Figure 5, are less to appear and impact the training process. Finally, to ameliorate the photometric loss term, the image \tilde{I}_t^{SF} is obtained by padding the $SF_{t \rightarrow s}$ at first, which is predicted at 416×128 , to original resolution (*e.g.*, 640×192), then using this flow to warp the full resolution I_s at I_t coordinates and finally extracting from this image the same crop as used before. This simple strategy allows to leverage on a complete image, since otherwise the cropped image would suffer from motion occlusions near boundaries. Moreover, we highlight that SD-OFNet is initialized to the OFNet weights, *i.e.* those found during the above described OFNet training based on the standard photometric loss, and then, when training SD-OFNet, only its weights are updated, *i.e.* OFNet is kept frozen.

3. Monocular Depth Estimation

In this section, we provide more insights on Ω Net performance concerning depth estimation, in particular by reporting comparison with state-of-the-art methods trained with *stronger* supervision, a more detailed analysis about the errors computed at different depth ranges and a reproducibility study about DSNet.

3.1. Comparison with more methods on the KITTI Eigen split

In this section, we report additional comparisons on the Eigen’s KITTI test split [11]. In particular, we compare Ω Net to state-of-the-art frameworks trained with stronger forms of supervision, *i.e.* stereo pairs, stereo videos or proxy labels. Differently from these approaches, we do not apply any post-processing step to further improve predictions. As highlighted in Table 3, we can notice how our method is comparable and, in most cases performs better, wrt other self-supervised depth-from-mono architectures trained on stereo pairs/stereo videos. Moreover, we point out that we outperform frameworks running online adaptation on the testing set [3, 7] on most metrics. Only semi-supervised methods at the bottom of the table [32, 27, 30] are in general more effective, because of the much stronger supervision from traditional stereo algorithms deployed during training.

3.2. Error at different depth ranges

In Table 3, we report more data supporting the claim that DSNet produces more accurate depth estimates at long distances with respect to other strategies such as [14] or even replacing our architecture with a much more complex one [34] based on a ResNet-50 backbone. We deeply looked into this and ascribe this finding to more complex models producing *over-smoothed* depth maps. In particular, in our experiments, we noticed that our shallow network tends to produce much sharper estimates compared to models having many more parameters. *Over-smoothing* produces better qualitative predictions and higher accuracy at short ranges, but it degrades depth accuracy at long distances, as we can observe in the table.

Method	M	S	V	P	A	I	CS	Lower is better				Higher is better		
								Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Zhou <i>et al.</i> [37]			✓				✓	0.198	1.836	6.565	0.275	0.718	0.901	0.960
Godard <i>et al.</i> [14]			✓			✓		0.115	0.903	4.863	0.193	0.877	0.959	0.981
Godard <i>et al.</i> [14]			✓					0.132	1.044	5.142	0.210	0.845	0.948	0.977
Godard <i>et al.</i> [14] (1024 × 320)			✓			✓		0.115	0.882	4.701	0.190	0.879	0.961	0.982
Zhou <i>et al.</i> [36]			✓			✓		0.121	0.837	4.945	0.197	0.853	0.955	0.982
Mahjourian <i>et al.</i> [19]			✓				✓	0.159	1.231	5.912	0.243	0.784	0.923	0.970
Yang <i>et al.</i> [33]			✓				✓	0.159	1.345	6.254	0.247	-	-	-
Wang <i>et al.</i> [28]			✓				✓	0.151	1.257	5.583	0.228	0.810	0.936	0.974
Bian <i>et al.</i> [2]			✓					0.128	1.047	5.234	0.208	0.846	0.947	0.970
Yin <i>et al.</i> [34]	✓		✓				✓	0.153	1.328	5.737	0.232	0.802	0.934	0.972
Zou <i>et al.</i> [38]	✓		✓				✓	0.146	1.182	5.215	0.213	0.818	0.943	0.978
Chen <i>et al.</i> [7]	✓		✓					0.135	1.070	5.230	0.210	0.841	0.948	0.980
Luo <i>et al.</i> [18]	✓		✓					0.141	1.029	5.350	0.216	0.816	0.941	0.976
Ranjan <i>et al.</i> [1]	✓		✓				✓	0.139	1.032	5.199	0.213	0.827	0.943	0.977
Xu <i>et al.</i> [31]			✓		✓			0.138	1.016	5.352	0.217	0.823	0.943	0.976
Casser <i>et al.</i> [3]			✓		✓			0.141	1.026	5.290	0.215	0.816	0.945	0.979
Gordon <i>et al.</i> [15]	✓		✓		✓			0.128	0.959	5.230	-	-	-	-
Ω Net(416 × 128)			✓		✓			0.134	0.893	5.137	0.208	0.829	0.946	0.979
Ω Net(640 × 192)			✓		✓			0.126	0.835	4.937	0.199	0.844	0.953	0.982
Ω Net(1024 × 320)			✓		✓			0.125	0.805	4.795	0.195	0.849	0.955	0.983
Ω Net(416 × 128)			✓		✓		✓	0.126	0.862	4.963	0.199	0.846	0.952	0.981
Ω Net(640 × 192)			✓		✓		✓	0.120	0.792	4.750	0.191	0.856	0.958	0.984
Ω Net(1024 × 320)			✓		✓		✓	0.118	0.748	4.608	0.186	0.865	0.961	0.985
Ω Net(768 × 384) †			✓		✓		✓	0.184	1.565	6.456	0.243	0.742	0.920	0.974
Casser <i>et al.</i> [3] (+ Online Ref.)			✓		✓			0.109	0.825	4.750	0.187	0.874	0.958	0.983
Chen <i>et al.</i> [7] (+ Online Ref.)			✓					0.099	0.796	4.743	0.186	0.884	0.955	0.979
Poggi <i>et al.</i> [22]		✓					✓	0.146	1.291	5.907	0.245	0.801	0.926	0.967
Poggi <i>et al.</i> [23]		✓					✓	0.111	0.849	4.822	0.202	0.865	0.952	0.978
Pillai <i>et al.</i> [21]		✓						0.112	0.875	4.958	0.207	0.852	0.947	0.977
Godard <i>et al.</i> [14]		✓	✓			✓		0.106	0.806	4.630	0.193	0.876	0.958	0.980
Godard <i>et al.</i> [13]		✓					✓	0.114	0.898	4.935	0.206	0.861	0.949	0.976
Zhang <i>et al.</i> [35]		✓	✓					0.135	1.132	5.585	0.229	0.820	0.933	0.971
Luo <i>et al.</i> [18]	✓	✓						0.127	0.936	5.008	0.209	0.841	0.946	0.979
Yang <i>et al.</i> [32]		✓	✓	✓				0.097	0.734	4.442	0.187	0.888	0.958	0.980
Watson <i>et al.</i> [30]		✓		✓		✓		0.096	0.710	4.393	0.185	0.890	0.962	0.981
Tosi <i>et al.</i> [27]		✓		✓			✓	0.096	0.673	4.351	0.184	0.890	0.961	0.981

Table 2. Quantitative evaluation on the Eigen test set of the KITTI dataset [12] for self-supervised monocular depth estimation methodologies. S: stereo pairs, V: video sequence, P: depth proxy labels, A: additional information, I: feature extractors pre-trained on ImageNet [9] or CS: Cityscapes [8]. †Trained on CS and tested on KITTI without any fine-tuning.

Method	Cap (m)	Abs Rel	Sq Rel	RMSE	RMSE log
Godard <i>et al.</i> [14]	0-8	0.059	0.062	0.503	0.082
Ours †	0-8	0.060	0.063	0.502	0.082
Ours	0-8	0.062	0.065	0.517	0.085
Godard <i>et al.</i> [14]	0-15	0.083	0.173	1.178	0.125
Ours †	0-15	0.083	0.168	1.148	0.122
Ours	0-15	0.084	0.169	1.156	0.124
Godard <i>et al.</i> [14]	0-30	0.111	0.470	2.561	0.172
Ours †	0-30	0.111	0.442	2.513	0.169
Ours	0-30	0.111	0.425	2.463	0.169
Godard <i>et al.</i> [14]	0-50	0.125	0.788	3.946	0.198
Ours †	0-50	0.127	0.762	4.020	0.199
Ours	0-50	0.124	0.702	3.836	0.195
Godard <i>et al.</i> [14]	0-80	0.132	1.044	5.142	0.210
Ours †	0-80	0.134	1.074	5.451	0.213
Ours	0-80	0.126	0.835	4.937	0.199

Table 3. Depth errors at different depth ranges. † indicates that our depth network has been replaced with the heavy-weight [34] backbone based on the ResNet50 architecture.

3.3. Reproducibility

We perform three independent training of our architecture to assess upon its reproducibility. Table 4 shows how our architecture produces the same results with negligible variance due to the randomness factors in training, *i.e.* initialization, data shuffle and augmentation.

Resolution	Abs Rel	Sq Rel	Lower is better		Higher is better		
			RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
640 × 192	0.120	0.792	4.750	0.191	0.856	0.958	0.984
640 × 192	0.122	0.799	4.749	0.191	0.856	0.958	0.984
640 × 192	0.121	0.795	4.755	0.192	0.855	0.957	0.983

Table 4. Three independent runs of our Ω Net(DSNet) result in slightly different models on the KITTI Eigen split.

4. Semantic Segmentation

In this section we report more detailed semantic segmentation results. Purposely, we use the following metrics:

1. **IoU**: Intersection over Union for pixel-wise segmentation calculated for each class or category, as defined in [8].
2. **mIoU_{class}**: mean IoU for the the 19 training classes used in CityScapes [8]: road, sidewalk, building, wall, fence, pole, traffic light, traffic sign, vegetation, terrain, sky, person, rider, car, truck, bus, train, motorcycle and bicycle.
3. **mIoU_{category}**, mean IoU considering the 7 macro-classes defined in CityScapes [8]: flat, construction, object, nature, sky, human, vehicles.
4. **Pixel Accuracy (Acc.)**: ratio between the correct and the total pixel predictions without considering any specific class or category.

4.1. Generalization across Datasets

In Tables 5 and 6 we validate with thorough experiments the motivation behind our better generalization across datasets compared to other state-of-the-art methods for real-time semantic segmentation. In this study we train on CityScapes and test on KITTI, reporting in Table 5 the IoU for the 19 classes, the mIoU_{class} and the pixel Acc. In Table 6 we report the IoU for the 7 categories and the mIoU_{category}. We refer with CS(S) methods trained on 2975 CityScapes images and with CS(P) methods trained on 22,973 proxy labels produced by [6]. To evaluate the performance of [16, 4] we used the official code and pre-trained weights available online. Our DSNet differs from other methods by three factors: 1) the architecture, 2) the training protocol exploiting proxy labels instead of ground truths and 3) the joint reasoning about geometry and semantics.

Regarding the tests on KITTI, our architecture trained only for semantic segmentation, namely Semantic Network or SNet, achieves good performance in Acc. but turns out worse than [4] for other metrics. On the other hand, it is worth to notice that training SNet with CS(P) allows our method to achieve a great performance boost in all metrics with respect to CS(P) (rows 8 vs 9). Finally, we can notice how DSNet achieves results comparable to SNet. This confirms the findings in [24], that joint reasoning about depth and semantics is more beneficial to the former task.

Method	Train	Test	road	sidewalk	building	wall	fence	pole	traffic light	traffic sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bicycle	mIoU _{class}	Acc.
DABNet [16]	CS(S)	CS	97.05	82.86	91.01	48.20	55.56	59.30	63.12	72.76	91.58	61.24	93.50	77.96	54.70	93.28	53.06	71.01	27.77	56.00	72.91	69.62	94.62
FCHardNet [4]	CS(S)	CS	97.39	84.40	92.31	53.83	62.90	64.28	68.21	78.06	91.85	59.82	94.91	80.81	60.55	94.85	72.70	82.15	76.45	59.97	75.49	76.37	95.35
ΩNet(SNet)	CS(S)	CS	93.69	65.66	83.46	23.57	20.57	40.11	35.32	47.77	86.62	44.22	89.94	56.00	23.00	84.98	17.22	1.22	0.00	17.11	52.82	46.49	89.56
ΩNet(SNet)	CS(P)	CS	95.97	77.23	87.96	38.37	42.62	47.82	48.15	60.44	89.73	54.97	92.62	65.87	36.96	90.57	25.19	0.06	0.00	25.53	61.06	54.80	92.45
ΩNet(DSNet)	CS(P)	CS	96.00	77.46	88.30	41.84	41.68	48.74	47.80	59.24	89.61	53.89	92.57	66.29	38.61	90.61	27.39	0.37	0.00	18.01	62.78	54.80	92.50
DABNet [16]	CS(S)	K	79.02	19.07	58.38	18.04	30.73	40.61	44.24	41.67	80.87	48.76	76.61	13.39	0.17	63.30	21.32	8.21	19.81	1.29	7.04	35.40	80.50
FCHardNet [4]	CS(S)	K	75.66	32.65	78.51	13.16	28.46	51.33	57.16	55.58	81.06	45.59	91.43	23.84	12.19	58.86	24.91	34.89	68.71	4.66	11.38	44.74	72.07
ΩNet(SNet)	CS(S)	K	83.31	33.39	66.57	12.15	20.18	44.20	37.76	32.35	84.46	58.79	88.70	24.66	13.55	76.09	12.62	2.09	0.10	1.15	12.64	37.09	84.94
ΩNet(SNet)	CS(P)	K	88.73	47.85	77.01	19.72	30.65	47.34	53.63	43.16	86.65	67.97	94.49	24.81	29.39	80.68	14.88	0.53	0.00	3.05	12.30	43.31	88.76
ΩNet(DSNet)	CS(P)	K	87.89	46.64	77.48	18.55	29.65	48.73	51.12	40.52	86.66	63.54	95.06	29.79	34.74	82.03	12.77	0.63	0.00	7.60	18.82	43.80	88.31

Table 5. IoU on 19 training classes, mIoU_{class} and pixel accuracy (Acc.) results of ΩNet against state of the art method training on CS and tested on CS or K. Better generalization from CS to K thanks to our proxy labels training protocol.

Method	Train	Test	flat	construction	object	nature	sky	human	vehicle	mIoU _{category}
DABNet [16]	CS(S)	CS	97.93	91.69	65.90	92.03	93.50	79.59	92.25	87.56
FCHardNet [4]	CS(S)	CS	98.19	92.55	70.77	92.27	94.91	82.31	93.54	89.22
ΩNet(SNet)	CS(S)	CS	96.34	84.29	44.37	86.85	89.94	60.13	83.77	77.96
ΩNet(SNet)	CS(P)	CS	97.40	88.80	53.61	90.19	92.62	69.08	88.47	82.88
ΩNet(DSNet)	CS(P)	CS	97.38	88.76	53.91	89.93	92.57	69.27	88.61	82.92
DABNet [16]	CS(S)	K	83.41	59.07	46.41	84.30	76.61	17.05	63.61	61.49
FCHardNet [4]	CS(S)	K	80.89	75.35	58.68	88.11	91.43	24.62	58.33	68.20
ΩNet(SNet)	CS(S)	K	87.93	63.92	45.79	85.47	88.70	31.02	69.95	67.54
ΩNet(SNet)	CS(P)	K	91.97	74.95	52.29	89.80	94.49	29.28	81.83	73.52
ΩNet(DSNet)	CS(P)	K	91.42	74.84	53.35	89.36	95.06	35.45	80.69	74.31

Table 6. IoU on 7 training categories and, mIoU_{category} results of ΩNet against state of the art method training on CS and tested on CS or K. Better generalization from CS to K thanks to our proxy labels training protocol.

4.2. Proxy Semantic Network

We evaluate the performance of the proxy semantic network. We employ DPC [6], pre-trained on CityScapes with the 2975 training ground truths. We report in Table 7 the testing results on the 500 and 200 images belonging to CityScapes validation set and the KITTI training datasets, respectively. Even though DPC [6] achieves impressive performance both on CityScapes as well as in generalizing to KITTI, it is a huge network unable to run in real-time (*i.e.*, it approximately delivers 3.5 fps on 768×384 images).

Method	Train	Test	mIoU _{class}	mIoU _{category}	Acc.
DPC[5] - Proxy	CS(S)	CS	80.22	90.73	95.99
DPC[5] - Proxy	CS(S)	K	58.75	81.30	90.21

Table 7. Semantic segmentation performances of the proxy semantic network [6] on CS and K datasets.

4.3. Priors Evaluation on KITTI

When we produce the priors used during training and, at prediction time, to create the M_t^{mot} , we split the 19 classes in static and potentially dynamic ones according to the following scheme:

1. **Static:** road, sidewalk, building, wall, fence, pole, traffic light, traffic sign, vegetation, terrain, sky
2. **Potentially dynamic:** person, rider, car, truck, bus, train

As among our objectives is to obtain a good motion segmentation mask, we are interested in evaluating the quality of our semantic segmentation predictions in terms of how they are amenable to producing good estimated priors according to the mapping defined above. We evaluate our DSNet trained on CityScapes+KITTI in the 200 KITTI images which provides semantic labels. We obtain a pixel accuracy of 98.50% while a 98.40% IoU for the static classes and a 80.99% for the potentially dynamic classes for a global 89.64% mIoU. It is worth noticing that, even though our segmentation is not able to perform a precise class segmentation, it yields excellent binary priors that turns out key to performance for motion segmentation.

5. Optical Flow Estimation

5.1. Comparison with more methods on the KITTI 2015 split

In Table 8 we include additional results from our main competitors to allow for a more comprehensive analysis. In particular, we report additional experiments from [1], in which the authors exploit a different combination of depth and optical flow networks, and from [29], that demonstrate how using stereo pairs at training time allows to obtain much better results on rigid regions. Nonetheless, it can be noticed that our network still outperforms existing monocular multi-task strategies by a large margin.

Method	Dataset	train			test
		Noc	All	F1	F1
Meister <i>et al.</i> [20] - C	SYN + K	-	8.80	28.94%	29.46%
Meister <i>et al.</i> [20] - CSS	SYN + K	-	8.10	23.27%	23.30%
Zou <i>et al.</i> [38]	SYN + K	-	8.98	26.0%	25.70%
Ranjan <i>et al.</i> [1] - DispResNet + PWC	SYN + K	-	5.66	20.93%	25.27%
Wang <i>et al.</i> [29] (Ego-motion) **	K	-	10.69	-	32.34%
Wang <i>et al.</i> [29] (Full) **	K	-	5.58	-	18.00%
Ren <i>et al.</i> [25]	K	-	16.79	36.00%	39.00%
Yin <i>et al.</i> [34]	K	8.05	10.81	-	-
Chen <i>et al.</i> [7] †	K	5.40	8.95	-	-
Chen <i>et al.</i> [7] (online) †	K	4.86	8.35	-	-
Ranjan <i>et al.</i> [1] - DispNet + FlowNetC	K	-	7.76	-	-
Ranjan <i>et al.</i> [1] - DispResNet + PWC	K	-	6.21	26.41%	-
Luo <i>et al.</i> [18]	K	5.84	-	21.56%	-
Luo <i>et al.</i> [18] *	K	5.43	-	20.61%	-
Ω Net (Ego-motion)	K	11.72	13.50	51.22%	-
Ω Net(OFFNet)	K	3.48	11.61	25.78%	-
Ω Net(SD-OFFNet)	K	3.29	5.39	20.00%	19.47%

Table 8. We report percentage of erroneous pixels (F1 score) and average end-point error over all pixels (All) and non-occluded pixels (Noc) on the KITTI 2015 flow dataset. We indicate with †feature extractors pre-trained on ImageNet, SYN as the SYNTHIA [26] dataset, CS for the Cityscapes dataset, multi-task methods *trained on stereo pair and ** using stereo at testing time.

6. Pose Estimation

We validate the performance of our framework on pose estimation on the KITTI odometry split, which provides ground-truth camera poses obtained with IMU/GPS readings for 11 driving sequences, indexed from 00 to 08 for training and 09-10 for testing purposes. As in [14], we have not changed our architecture for this specific task but simply trained it from scratch on new training sequences without known intrinsic parameters. We compare our model with learned camera intrinsic parameters with several monocular self-supervised methods on the two sequences of KITTI odometry test split. All of the results, summarized in 9, are evaluated by optimizing the scaling factor to align with the ground-truth to address the inherent scale ambiguity.

Method	Frames	Sequence 09	Sequence 10
Zhou <i>et al.</i> [37]	5	0.021 ± 0.017	0.020 ± 0.015
Ranjan <i>et al.</i> [1]	5	0.012 ± 0.007	0.012 ± 0.008
Yin <i>et al.</i> [34]	5	0.012 ± 0.007	0.012 ± 0.009
ORB-Slam	3	0.014 ± 0.008	0.012 ± 0.011
Casser <i>et al.</i> [3]	3	0.011 ± 0.006	0.011 ± 0.010
Zou <i>et al.</i> [3]	3	0.017 ± 0.007	0.015 ± 0.009
Luo <i>et al.</i> [18]	3	0.013 ± 0.007	0.012 ± 0.008
Godard <i>et al.</i> [14]	2	0.017 ± 0.008	0.015 ± 0.010
Ours †	2	0.020 ± 0.013	0.017 ± 0.011

Table 9. Absolute Trajectory Error (ATE) of pose estimation evaluated on the KITTI odometry split sequences 09-10. †indicates strategies trained with unknown camera intrinsics.

7. Motion Segmentation

7.1. Threshold analysis

In Figure 1, we present an ablation study dealing with the motion segmentation task. In the main paper, to be consistent with other methodologies, we set the threshold τ used for the evaluation to 0.5. However, we point out that a careful tuning of such threshold can improve the overall motion segmentation accuracy. In particular, we can notice how the best configuration for our predictions is obtained using a larger threshold. Indeed, we found out that the best trade-off between the mean accuracy and the mean IoU is achieved by setting the threshold value to 0.7 (in this case the Mean Acc is 0.91 while Mean IoU is 0.77).

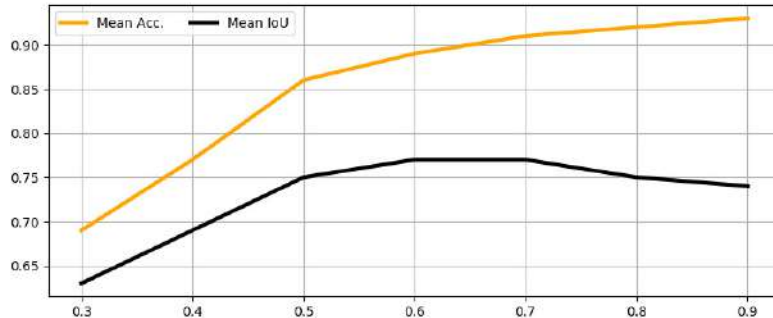


Figure 1. Mean Acc. and mIoU varying the threshold used to compute the motion segmentation M_t^{mot} .

7.2. Evaluation for KITTI only on Cars

We conduct an additional study to evaluate our motion segmentation masks only on pixels belonging to Cars, as proposed in [1]. In Table 10 we evaluate the IoU for static and dynamic cars yielded by Ω Net and [1] on the 200 KITTI images endowed with ground truth for the motion segmentation task. We notice that our M_t^{mot} outperforms [1] in all metrics (rows 1 vs 2 and 3) for all thresholds. Moreover, we point out that in this test configuration the contribution given to the motion segmentation by our estimated semantics is almost negligible as car regions are already extracted by using KITTI ground truths. Therefore, we test also our motion probability P_t alone, showing that it is superior to [1] even without the help provided by the estimated semantics.

Method	Threshold	Overall	Static Cars	Moving Cars
Ranjan [1]	-	56.94	55.77	58.11
Ω Net M_t^{mot}	0.5	63.98	64.16	63.79
Ω Net M_t^{mot}	0.7	63.97	64.15	63.79
Ω Net P_t	0.5	63.67	62.58	64.77
Ω Net P_t	0.7	62.66	58.42	66.89

Table 10. Motion Segmentation Results. IoU scores on KITTI 2015 training dataset images computed only over car pixels.

8. Runtime

In this section we report additional runtime results on the three different devices used in the main paper, that is: an NVIDIA Titan Xp GPU, an Intel i7-7700K CPU and an NVIDIA Jetson TX2 GPU. In Table 11, we show further timings by varying the input image resolution of our architecture. It can be noticed how Ω Net achieves real-time results (*i.e.* 27.9) on the Titan Xp GPU even with the largest image size 1024×320 , reaching about 2 FPS on the Jetson Tx2 embedded device with the same input configuration.

	416 × 128						640 × 192					1024 × 320				
	W	D	DS	OF	Cam	O	D	DS	OF	Cam	O	D	DS	OF	Cam	O
Jetson TX2	15	20.2	17.9	8.9	54.1	7.1	12.5	10.3	6.5	49.2	4.5	6.4	5.3	3.2	26.31	2.0
i7-7700K	91	10.9	9.1	11.0	60.1	5.5	5.0	4.2	4.9	31.4	2.4	1.9	1.6	1.8	13.2	0.9
Titan XP	250	250.7	212.4	152.6	550.7	90.5	170.2	134.1	94.1	446.7	57.4	86.0	64.5	44.5	251.0	27.9

Table 11. Runtime analysis on different hardware devices. For each device we report the power consumption in Watt and the FPS by varying input resolution. D: Depth, S: Semantic, OF: Optical Flow, Cam: camera pose, O: Overall architecture.

9. Qualitative results

In Figures 2, 3, 4, 5, 6, 7, 8, we provide qualitative results of our architecture on the standard datasets used in the main paper, such as KITTI and CityScapes. We refer the reader to the captions for description and comments related to each example.

9.1. Results on a YouTube Video

Furthermore, to prove that our network can be trained on unconstrained monocular sequences with unknown camera parameters and without semantic ground-truth labels, we downloaded from YouTube an online video captured by a moving camera consisting of 130K images depicting an urban scenario. Then, we generated proxy semantic labels using [6] and trained Ω Net(DSNet) to learn depth, pose, semantics and camera intrinsics. Figure 9, show qualitative results yielded by Ω Net on this unconstrained monocular video.

References

- [1] Ranjan Anurag, Varun Jampani, Kihwan Kim, Deqing Sun, Jonas Wulff, and Michael J. Black. Competitive collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 4, 7, 8
- [2] Jia-Wang Bian, Zhichao Li, Naiyan Wang, Huangying Zhan, Chunhua Shen, Ming-Ming Cheng, and Ian Reid. Unsupervised scale-consistent depth and ego-motion learning from monocular video. In *Thirty-third Conference on Neural Information Processing Systems (NeurIPS)*, 2019. 4
- [3] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In *Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)*, 2019. 3, 4, 7
- [4] Ping Chao, Chao-Yang Kao, Yu-Shan Ruan, Chien-Hsiang Huang, and Youn-Long Lin. Hardnet: A low memory traffic network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3552–3561, 2019. 5, 6
- [5] Liang-Chieh Chen, Maxwell Collins, Yukun Zhu, George Papandreou, Barret Zoph, Florian Schroff, Hartwig Adam, and Jon Shlens. Searching for efficient multi-scale architectures for dense image prediction. In *Advances in Neural Information Processing Systems*, pages 8699–8710, 2018. 6
- [6] Liang-Chieh Chen, Maxwell D. Collins, Yukun Zhu, George Papandreou, Barret Zoph, Florian Schroff, Hartwig Adam, and Jonathon Shlens. Searching for efficient multi-scale architectures for dense image prediction. In *NIPS*, 2018. 1, 5, 6, 9
- [7] Yuhua Chen, Cordelia Schmid, and Cristian Sminchisescu. Self-supervised learning with geometric constraints in monocular video: Connecting flow, depth, and camera. In *ICCV*, 2019. 3, 4, 7
- [8] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 4, 5
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 4
- [10] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *ICCV*, 2015. 12
- [11] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE international conference on computer vision*, pages 2650–2658, 2015. 3
- [12] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 4
- [13] Clément Godard, Oisín Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, 2017. 1, 4
- [14] Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *ICCV*, 2019. 1, 2, 3, 4, 5, 7
- [15] Ariel Gordon, Hanhan Li, Rico Jonschkowski, and Anelia Angelova. Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. In *ICCV*, 2019. 4

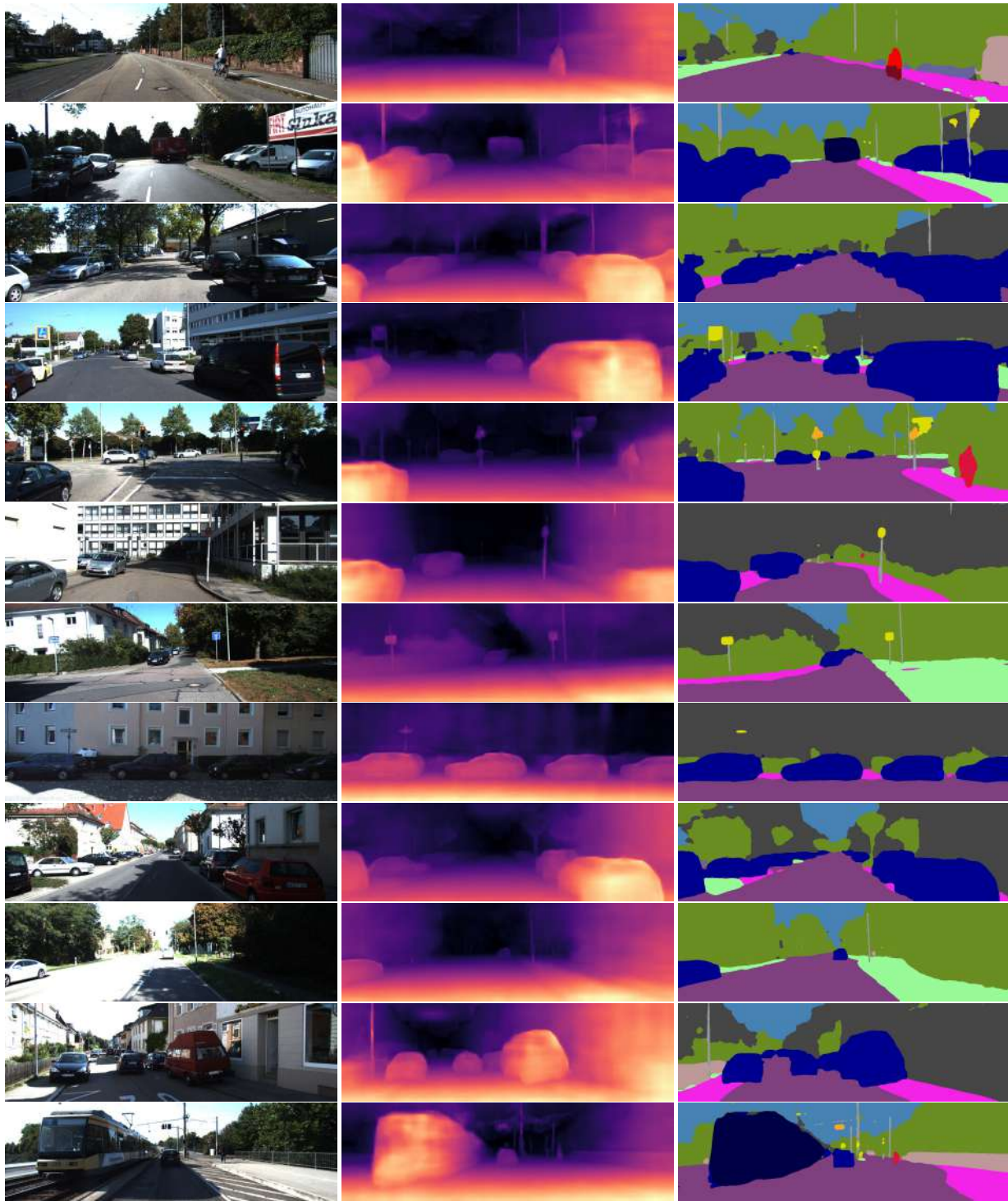


Figure 2. Qualitative results of our Ω Net(DSNet) on the KITTI Eigen split. From left to right we show image, the predicted single-image depth map and the predicted semantic segmentation of the scene.

- [16] Gen Li and Joongkyu Kim. Dabnet: Depth-wise asymmetric bottleneck for real-time semantic segmentation. In *British Machine Vision Conference*, 2019. 5, 6
- [17] Pengpeng Liu, Michael R. Lyu, Irwin King, and Jia Xu. Selfflow: Self-supervised learning of optical flow. In *CVPR*, 2019. 1
- [18] Chenxu Luo, Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, Ram Nevatia, and Alan Yuille. Every pixel counts++: Joint learning

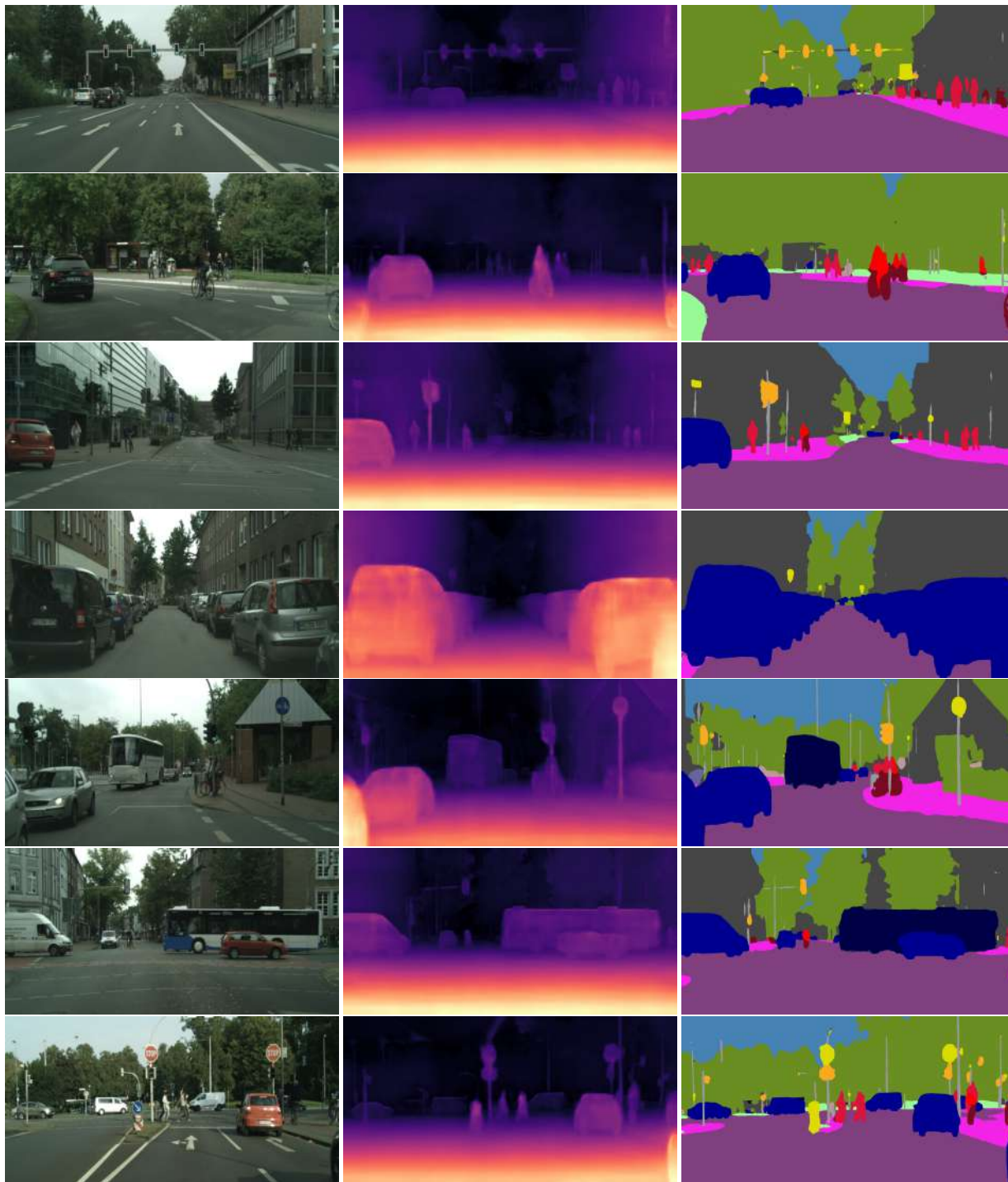


Figure 3. Qualitative results of our network Ω Net(DSNet) on the Cityscapes dataset. From left to right, the input image, single-view depth and the semantic prediction of our model.

of geometry and motion with 3d holistic understanding. *PAMI*, 2019. 4, 7

- [19] Reza Mahjourian, Martin Wicke, and Anelia 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)*, 2018. 4
- [20] Simon Meister, Junhwa Hur, and Stefan Roth. UnFlow: Unsupervised learning of optical flow with a bidirectional census loss. In *AAAI*, New Orleans, Louisiana, Feb. 2018. 7, 12

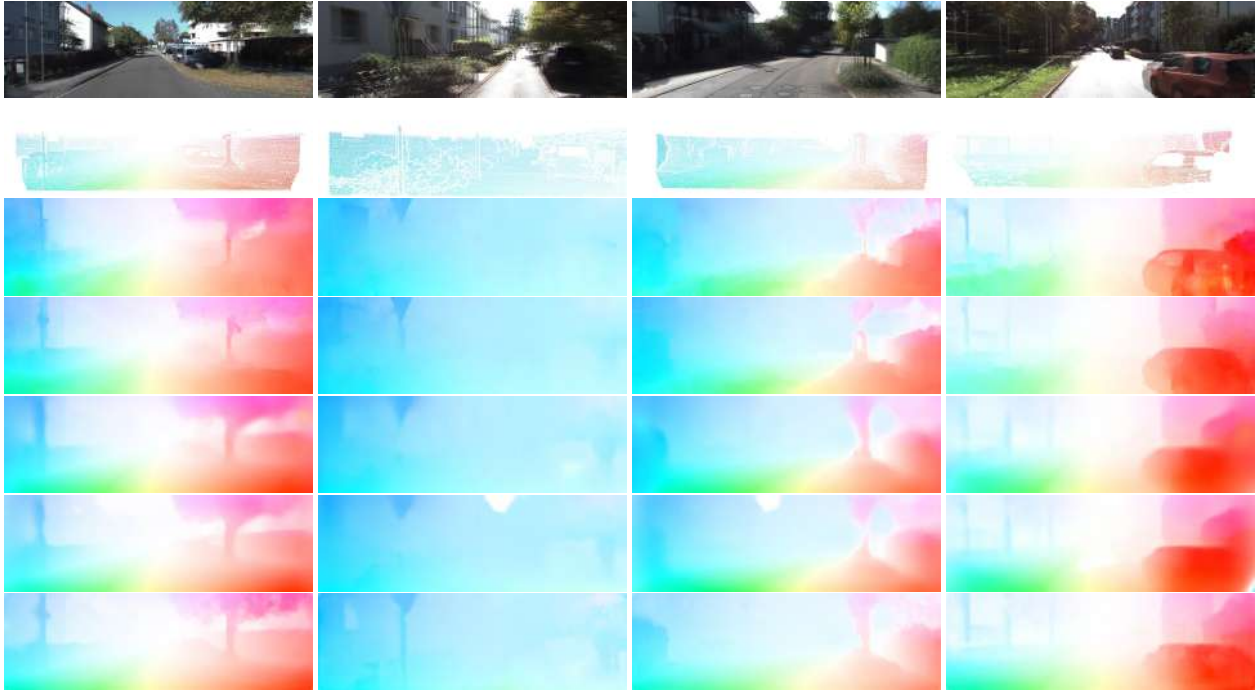


Figure 4. Qualitative results on optical flow estimation using the KITTI 2015 dataset. From top to bottom we show image, ground-truth labels, FlowNetS [10], FlowNetC [10], Unflow [20], DF-Net [38] and our Ω Net(SD-OFNet).

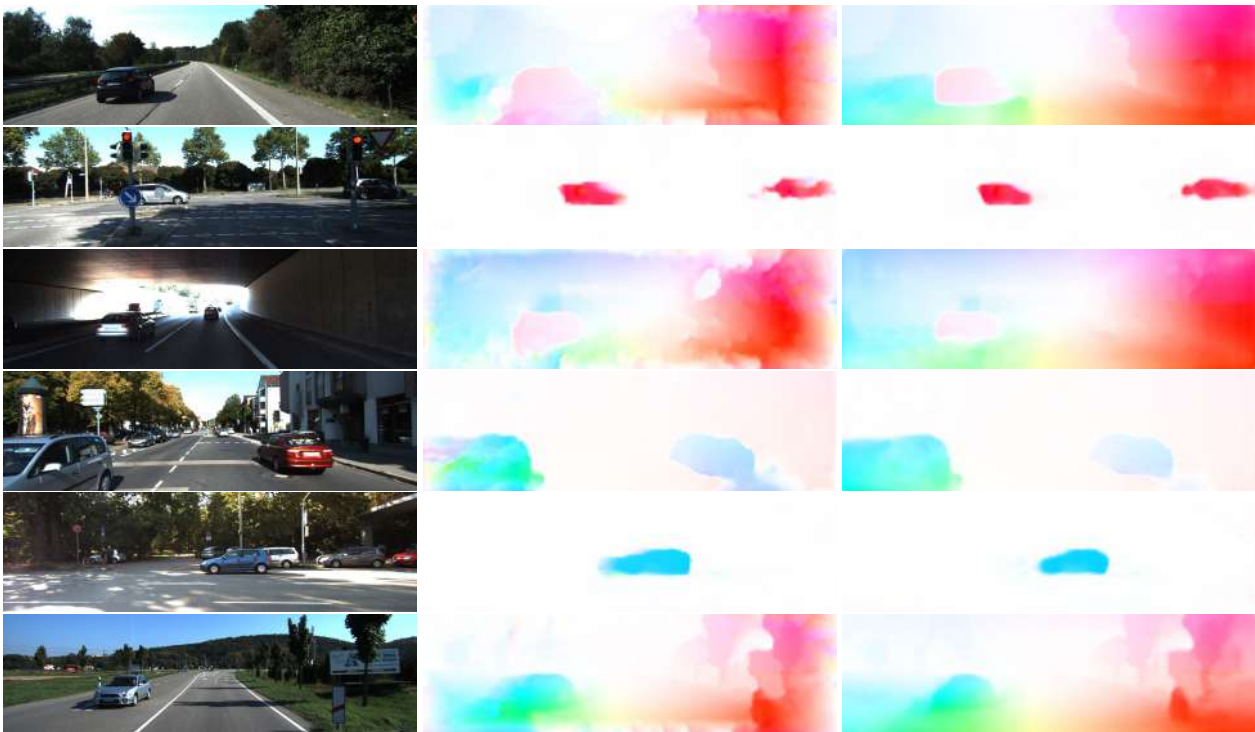


Figure 5. Qualitative comparison between the initial optical flow network Ω Net(OFNet) and the self-distilled one Ω Net(SD-OFNet) obtained with our strategy on the KITTI 2015 dataset. From left to right we show image, Ω Net(OFNet) and Ω Net(SD-OFNet) results respectively. It can be noticed that the proposed self-distilled paradigm greatly alleviates motion boundaries occlusions and improves details in the final optical flow predictions.

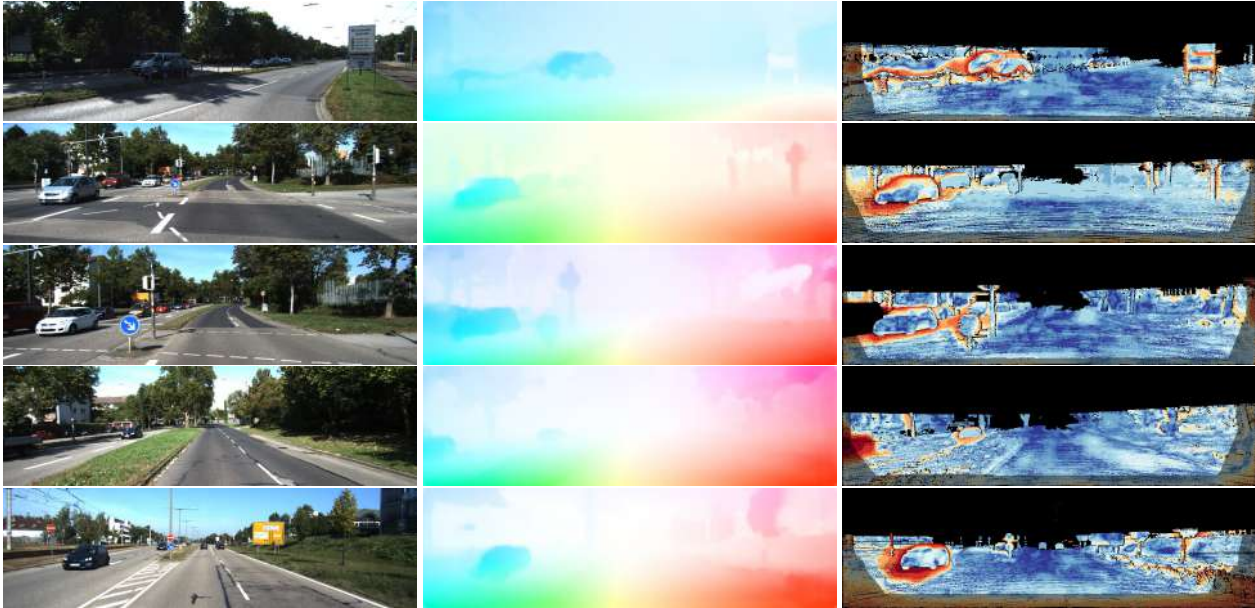


Figure 6. Visualization of the flow error map of our optical flow network on the KITTI 2015 testing benchmark. Larger errors are encoded in red, while blue pixels represents good optical flow estimates with respect to the ground-truth.

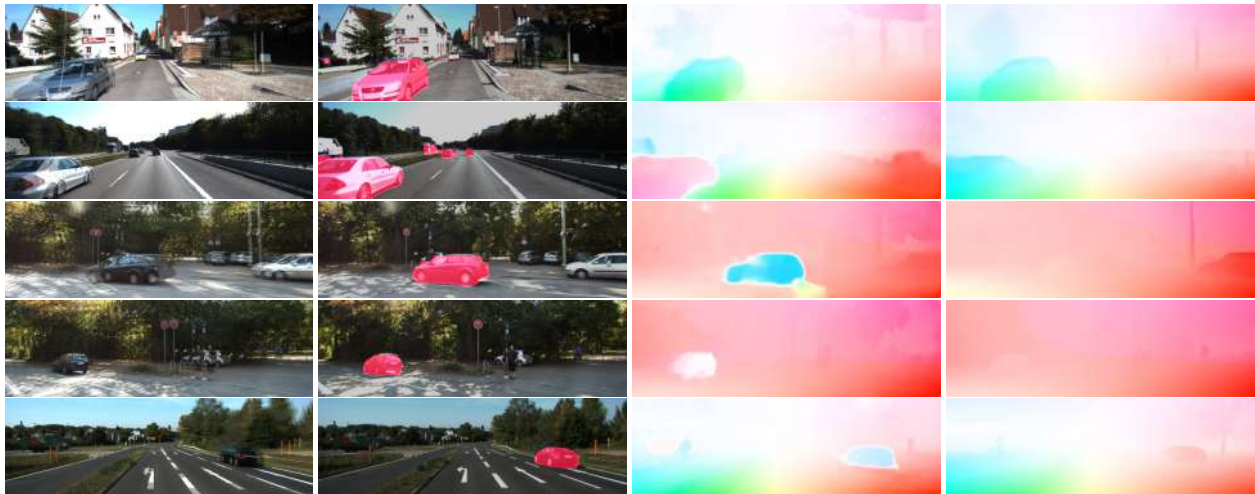


Figure 7. Qualitatives on motion segmentation. From left to right, the input images, the motion objects detected in the scene by our method (highlighted in red), the optical flow and the rigid flow.

- [21] Sudeep Pillai, Rares Ambrus, and Adrien Gaidon. Superdepth: Self-supervised, super-resolved monocular depth estimation. In *ICRA*, 2019. 4
- [22] Matteo Poggi, Filippo Aleotti, Fabio Tosi, and Stefano Mattoccia. Towards real-time unsupervised monocular depth estimation on CPU. In *IEEE/RSJ Conference on Intelligent Robots and Systems (IROS)*, 2018. 4
- [23] Matteo Poggi, Fabio Tosi, and Stefano Mattoccia. Learning monocular depth estimation with unsupervised trinocular assumptions. In *6th International Conference on 3D Vision (3DV)*, 2018. 4
- [24] Pierluigi Zama Ramirez, Matteo Poggi, Fabio Tosi, Stefano Mattoccia, and Luigi Di Stefano. Geometry meets semantics for semi-supervised monocular depth estimation. In *Asian Conference on Computer Vision*, pages 298–313. Springer, 2018. 2, 5
- [25] Zhe Ren, Junchi Yan, Bingbing Ni, Bin Liu, Xiaokang Bin, and Hongyuan Zha. Unsupervised deep learning for optical flow estimation. In *Thirty-First AAAI Conference on Artificial Intelligence (AAAI-17)*, 2017. 7
- [26] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio Lopez. The SYNTHIA Dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *CVPR*, 2016. 7

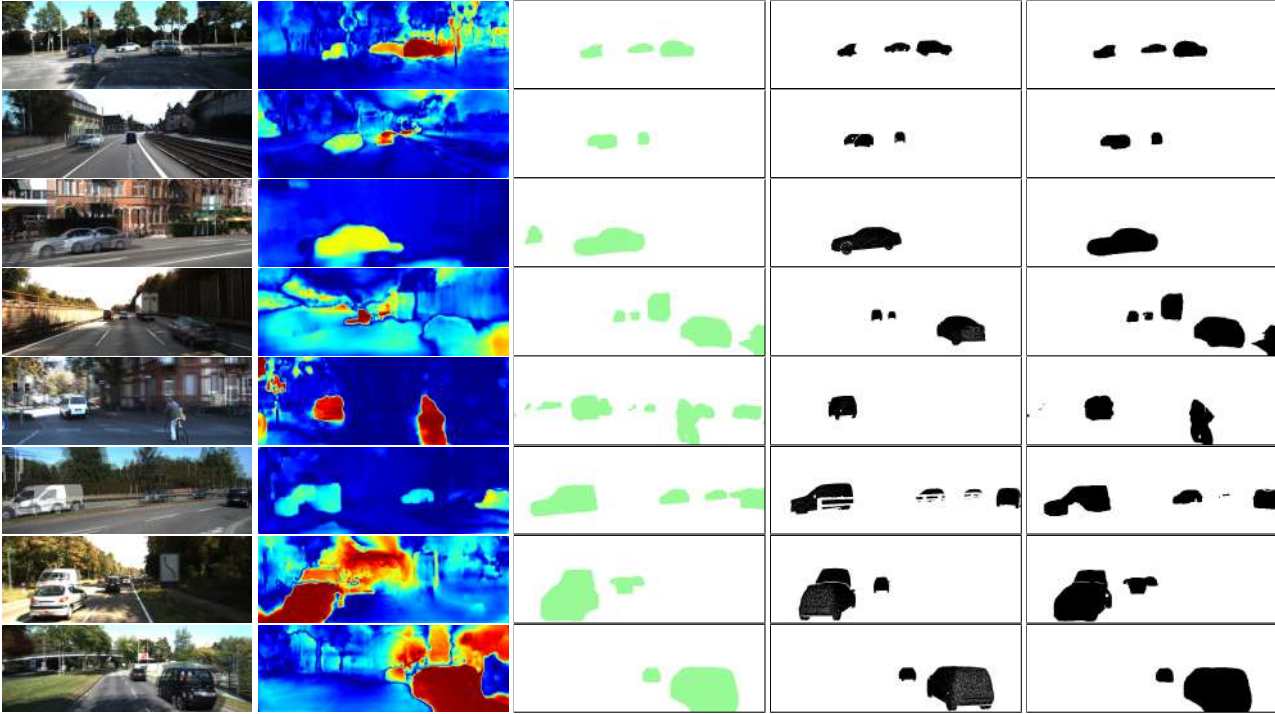


Figure 8. Motion segmentation results on the KITTI 2015 dataset. From left to right, we show the monocular sequence, the outcome of the proposed motion probability strategy (high probability of motion is encoded in red), semantic priors extracted from ours semantic predictions, ground-truth motion masks and ours estimated motion segmentation masks.

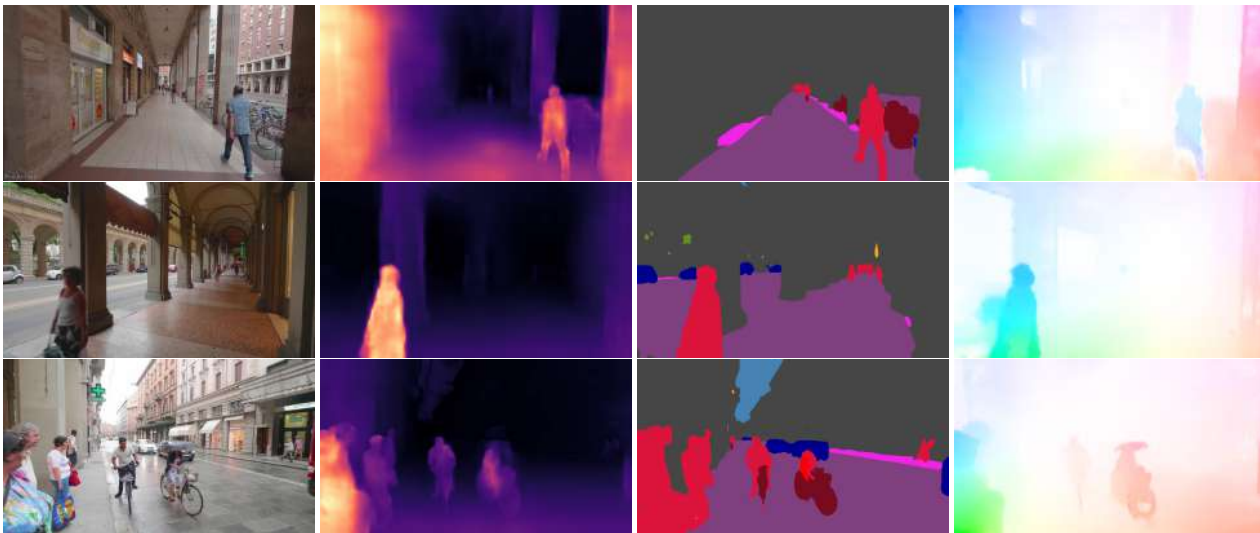


Figure 9. Qualitative results of Ω Net on a raw YouTube video. From left to right, we show the input images of a monocular sequence, the single-view depth and semantic predictions and, finally, the optical flow estimate.

[27] Fabio Tosi, Filippo Aleotti, Matteo Poggi, and Stefano Mattoccia. Learning monocular depth estimation infusing traditional stereo knowledge. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3, 4

[28] Chaoyang Wang, Jose Miguel Buenaposada, Rui Zhu, and Simon Lucey. Learning depth from monocular videos using direct methods. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 4

[29] Yang Wang, Peng Wang, Zhenheng Yang, Chenxu Luo, Yi Yang, and Wei Xu. Unos: Unified unsupervised optical-flow and stereo-depth estimation by watching videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8071–8081, 2019. 7

- [30] Jamie Watson, Michael Firman, Gabriel J Brostow, and Daniyar Turmukhambetov. Self-supervised monocular depth hints. In *ICCV*, 2019. 3, 4
- [31] Haofei Xu, Jianmin Zheng, Jianfei Cai, and Juyong Zhang. Region deformer networks for unsupervised depth estimation from unconstrained monocular videos. In *IJCAI*, 2019. 4
- [32] Nan Yang, Rui Wang, Jörg Stückler, and Daniel 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. 3, 4
- [33] Zhenheng Yang, Peng Wang, Wang Yang, Wei Xu, and Nevatia Ram. Lego: Learning edge with geometry all at once by watching videos. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 4
- [34] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3, 4, 5, 7
- [35] Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian 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)*, 2018. 4
- [36] Junsheng Zhou, Yuwang Wang, Naiyan Wang, and Wenjun Zeng. Unsupervised high-resolution depth learning from videos with dual networks. In *Inter. Conf. on Computer Vision*. IEEE, IEEE, October 2019. 4
- [37] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G. Lowe. Unsupervised learning of depth and ego-motion from video. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 4, 7
- [38] Yuliang Zou, Zelun Luo, and Jia-Bin Huang. Df-net: Unsupervised joint learning of depth and flow using cross-task consistency. In *European Conference on Computer Vision (ECCV)*, 2018. 4, 7, 12