

Learning Independent Object Motion from Unlabelled Stereoscopic Videos

Zhe Cao
UC Berkeley

zhecao@berkeley.edu

Abhishek Kar
Fyusion Inc.

akar@fyusion.com

Christian Häne
Google

chaene@google.com

Jitendra Malik
UC Berkeley

malik@eecs.berkeley.edu

Abstract

We present a system for learning motion maps of independently moving objects from stereo videos. The only annotations used in our system are 2D object bounding boxes which introduce the notion of objects in our system. Unlike prior learning based approaches which have focused on predicting dense optical flow fields and/or depth maps for images, we propose to predict instance specific 3D scene flow maps and instance masks from which we derive a factored 3D motion map for each object instance. Our network takes the 3D geometry of the problem into account which allows it to correlate the input images and distinguish moving objects from static ones. We present experiments evaluating the accuracy of our 3D flow vectors, as well as depth maps and projected 2D optical flow where our jointly learned system outperforms earlier approaches trained for each task independently.

1. Introduction

Consider the crowded road scene in Figure 1, what information do we as humans use to navigate effectively in this environment? We need to have an understanding of the structure of the environment, i.e. how far other elements in the scene (cars, bikes, people, trees) are from us. Moreover, we also require knowledge of the speed and direction in which other agents in the environment are moving relative to us. Such a representation, in conjunction with our ego-motion, enables us to produce a hypothesis of the environment state in the near future and ultimately allows us to plan our next actions.

In order to gather this information, humans use stereo-motion, i.e. a stream of images captured with our two eyes as we move through the environment. In this work, we develop a computational system that aims to produce such a factored scene representation of 3D structure and motion from a binocular video stream. Specifically, we propose to predict the 3D object motion of each moving object (represented by 3D scene flow) in addition to a detailed depth map of the scene from a stereo image sequence. This task and



Figure 1. Object motion predicted by our system. Trained with raw stereo motion sequences in a self-supervised manner, our model learns to predict object motion together with the scene depth using sequence of stereo images and object proposals as input. The speed and moving direction of each moving object is derived from our scene flow prediction.

its variants have been tackled in supervised settings which require labels such as dense depth maps and motion annotations that are prohibitively expensive to collect or alternatively obtained from synthetic datasets [4, 5, 17, 21, 27]. We present a system that learns to predict these quantities using only unlabelled stereo videos, thus making it applicable at scale. In addition to producing pixel-wise depth and scene flow maps, our network is aware of the notion of independent objects. This allows us to produce a rich factored 3D representation of the environment where we can measure velocities of independent objects in addition to their 3D positions and extents in the scene. The only labels used by our system are those introduced by off-the-shelf object detectors which are very cheap to acquire at scale.

Prior work in this domain has focused on certain sub-problems such as learning depth or optical flow prediction without explicit labels [49, 12, 8]. In Section 5, we demonstrate that by jointly learning the full problem of depth and scene flow prediction, we outperform these methods for each of these sub-problems as well. The key contributions of our work are as follows: (1) formulating a learning objective which works with the limited amount of supervision that can be gathered in a real world scenario (object bounding box annotations), (2) factoring the scene repre-

sentation into independently moving objects for predicting dense depth and 3D scene flow and (3) designing a network architecture that encodes the underlying 3D structure of the problem by operating on plane sweep volumes.

The sections in this paper are organized as follows. Section 2 discusses prior work on inferring scene structure and motion. Section 3 presents our technical approach for inferring scene flow from stereo motion - loss functions, object-centric prediction and priors. In Section 4, we describe our network architecture designed for geometric matching and 3D reasoning in plane sweep volumes. Section 5 details our experiments on the KITTI dataset [29] with extensive evaluation of our depth and scene flow prediction.

2. Related work

In our work we recover scene geometry and object motion jointly while traditionally these problems have been solved independently. The geometry of a scene is reconstructed by first recovering the relative camera pose between two or more images taken from different viewpoints using Structure-from-Motion (SfM) techniques [25, 14]. Subsequently, with dense matching and triangulation a dense 3D model of the scene is recovered [31]. The underlying assumption within the aforementioned methods is that the scene is static, i.e. does not contain moving objects. The case for independently moving objects has been studied in a purely geometric setting [3]. The key difficulties are degenerate configurations and outliers in point correspondences [30]. Therefore additional priors are used - a common example is objects moving on a ground plane [50]. Similarly, estimating the shape of non-rigid objects is ambiguous and hence using additional constraints such as maximizing the rigidity of the shape [41] or representing the non rigid shape as linear combination of base shapes [2] have been proposed. When reconstructing videos captured in unconstrained environments additional difficulties such as incomplete feature tracks and bleeding into the background have to be handled [6]. Our proposed approach is trained on real world data which makes it robust to appearance variations and suitable priors are directly learned from data.

Vedula *et al.* [42] introduced the problem of 3D scene flow estimation, where for each point a 3D motion vector between time t and $t + 1$ is computed. Different variants are considered depending on the amount of 3D structure that is given as input. A common variant is to consider a stream of binocular image pairs of a moving camera as input [16, 46, 44, 29, 38], and give a depth and 3D scene flow as output. This is often referred to as the stereo scene flow estimation problem. Similarly RGBD scene flow considers a stream of RGBD (color and depth) images as input [18].

Recently learning-based approaches, especially convolutional neural networks have been applied for single view depth prediction [23, 4], optical flow [5], stereo matching

and scene flow [27]. These learning systems are trained using ground truth geometry and/or flow data. In practice such data is only available for synthetic data in a large scale. A natural way to complement the limited amount of ground truth data is using weaker supervision. For the aforementioned problems, loss functions which are purely based on images and rely on photometric consistency as learning objective have been proposed [8, 51, 12, 39, 43]. They essentially utilize a classical non-learned system [7] within the loss function. A few recent works [49, 52, 47, 26, 33] use such a self-supervised approach to predict optical flow and depth. To our knowledge our work is the first network that learns to directly predict object specific 3D scene flow without relying on pixel-wise flow or depth annotations.

Another key difference of our work from prior works that predict depth and optical flow is that they predict depth based on a single image. This limits their performance as demonstrated in our results. Geometric reasoning can be included into the network architecture as demonstrated in [21, 20, 19, 48]. We extend these ideas to full 3D scene flow estimation while also operating at the level of object instances allowing us to produce rich factored geometry and motion representation of the scene.

3. Scene Flow from Stereo Motion

Figure 2 illustrates our system. A stream of calibrated binocular stereo image pairs $\mathcal{I} = \{I_1^l, I_1^r, \dots, I_n^l, I_n^r\}$ captured from times 1 to n is given as input. The most common case we are investigating is $n = 2$, i.e. two binocular frames at time t and $t + 1$. The intrinsic camera calibration K is assumed to be known. The camera poses of the left camera at each time instant are denoted by $\mathcal{T} = \{T_1, \dots, T_n\}$ and are precomputed using visual SLAM [10]. For any time instant t , we also have a set of j 2D bounding box detections $\mathcal{B} = \{B^1, \dots, B^j\}$ on the left image I_t^l predicted by an off-the-shelf object detector. The task is to compute the following quantities for the reference frame - a dense depth map D , a set of dense 3D flow fields $\mathcal{F} = \{F^1, \dots, F^j\}$ that describe the motion between t and $t + 1$ and a set of instance masks $\mathcal{M} = \{M^1, \dots, M^j\}$ for each moving object. From these instance-level predictions, we can compose the full scene flow map F by assigning a 3D scene flow vector to each image pixel in the full image.

We design our system as a convolutional neural network (CNN) which learns to predict all quantities jointly and train the network in a self-supervised manner. The supervision comes from the consistency between synthesized images and input images at different time instants and from different camera viewpoints. The basic principle is that given the predictions of the scene flow F and depth D in a frame I_{ref} , we can use the precomputed ego-motion to warp another image I into the reference view. This process generates a synthesized image which we call \hat{I} . We then define

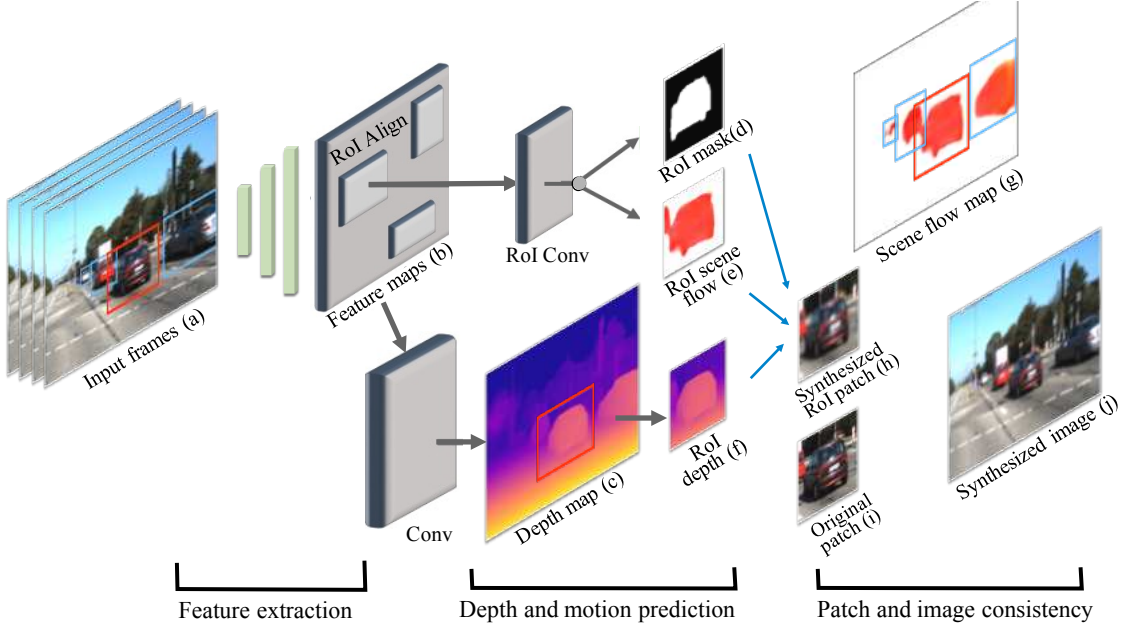


Figure 2. Our pipeline for learning depth and object motion. Using a stereo motion sequence as input, our system predicts a depth map (c), instance mask (d) and 3D scene flow (e) for each independent moving object in a single forward pass. Using the instance mask and scene flow, we compose a full scene flow map (g). For each region of interest (RoI), we synthesize a patch (h) based on the RoI camera intrinsics, RoI depth (f), 3D scene flow (e) and instance mask (d) as explained in Section 3.2. We use the synthesized patch (h) and original patch (i) from the input image to enforce consistency losses to supervise the RoI prediction. We use stereo reprojection to supervise the depth prediction. Finally, we use the full map scene flow and depth to synthesize a image (j) for computing the consistency loss.

our learning objective as the similarity between the captured images I_{ref} and the synthesized images \hat{I} . The above principle is then applied to each region of interest (RoI) independently followed by an assembly procedure for full image scene flow. This allows us to produce a factored representation of the environment into static and dynamic objects with high-quality estimates of instance masks, depth and motion.

3.1. Disentangling Camera and Object Motion

The motion in a dynamic scene captured by a moving camera can be decomposed into two elements - the motion of static background resulting from the camera motion and the motion of independently moving objects in the scene. A common way to represent the scene motion is 2D optical flow. However, this representation confounds the camera and object motion. We model the motion of the static background using the 3D structure represented as a depth map and the camera motion. Dynamic objects are modelled with full 3D scene flow. To this end, we utilize 2D object detections in the form of bounding boxes and reason about the 3D motion of each object independently.

3.2. Supervising Scene Flow by View Synthesis

The key supervision for the scene flow prediction comes from the photometric consistency of multiple views of the same scene. The process is illustrated in Figure 3. Our

network predicts a depth map D and a scene flow map F for the reference view I_{ref} . Using a different image I we can use the predictions to warp I into the reference view and generate a synthesized image \hat{I} . We then minimize the photometric difference between I_{ref} and \hat{I} given as

$$\mathcal{L}_{\text{photo}} = \alpha \frac{1 - \text{SSIM}(I_{\text{ref}}, \hat{I})}{2} + (1 - \alpha) \|I_{\text{ref}} - \hat{I}\|_1 \quad (1)$$

where SSIM denotes the structural similarity index [45] and α denotes a weighting parameter.

We denote the homogeneous coordinates of pixel p as $h(p)$. A pixel p from the reference frame is transformed to a pixel \hat{p} within a frame I

$$h(\hat{p}) = K T_{\text{rel}}(D(p) K^{-1} h(p) + F(p)) \quad (2)$$

with T_{rel} the relative transformation from reference frame to I . This allows us to do a reverse warp using bilinear interpolation, keeping the formulation differentiable.

Using the photometric consistency alone is insufficient for supervising the 3D flow prediction. The reason is that along a viewing ray multiple photo consistent solutions are possible, as shown in Figure 3. Therefore we use an additional geometric loss leveraging depth consistency which further constrains the flow. The idea is that the flow in z -direction, sometimes also called disparity difference has to

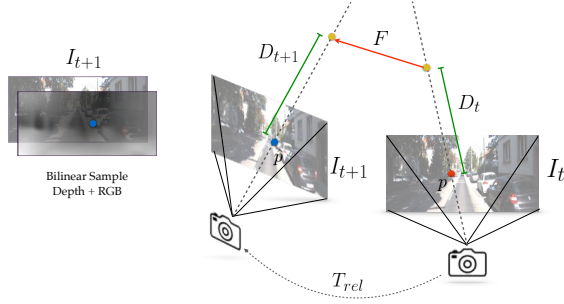


Figure 3. Illustration of our image reprojection process. A pixel p from image I_t is unprojected using its predicted depth and subsequently transformed to the frame of I_{t+1} using the predicted flow F and the camera transform T_{rel} . The photometric consistency loss is derived from the photometric difference between I_t and $\hat{I}_{t+1 \rightarrow t}$ where $\hat{I}_{t+1 \rightarrow t}$ is created by warping I_{t+1} into I_t . The geometric consistency loss is computed by comparing the difference between depth maps warped in the above manner and having them consistent with the z -dimension of the predicted flow F . Note that using only photometric consistency would not resolve the ambiguity in the z direction of the flow.

agree with the depth maps predicted for the two time instants t and $t + 1$. In order to utilize this loss function a depth map for both time instants needs to be predicted and the warping is applied to the depth map.

Analogous to the photometric consistency, the geometric consistency is defined by comparing the predicted depth values of the warped image and reference image,

$$\mathcal{L}_{\text{geo}} = \left\| D_{\text{ref}} - \hat{D} + F_z \right\|_1 \quad (3)$$

where D_{ref} refers to the predicted depth at time t and \hat{D} is the predicted depth at time $t + 1$ warped back to time t , F_z is the z -dimension of the predicted scene flow.

3.3. Object-centric Scene Flow Prediction

Image based consistency losses are typically applied by warping the whole image and then computing the consistency over the whole image - examples for optical flow prediction can be found in [49, 52]. For 3D scene flow this is not an ideal choice due to the sparsity of non-zero flow vectors. Compared to the static background, moving objects constitute only a small fraction of the image pixels. This unbalanced moving/static pixel distribution makes naively learning full image flow hard and ends up in zero flow predictions even on moving objects. To make the network focus on predicting the correct flow on moving objects and provide a more balanced supervision, we therefore use object bounding box detections obtained from a state-of-the-art 2D object detection system [24]. It is important to note that the object detection does not actually tell us if the object is moving or not. This information is learned by our network using our view synthesis based loss functions.

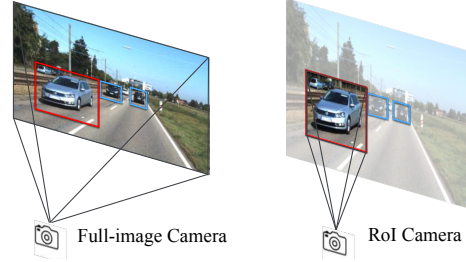


Figure 4. Illustration of image rescale and crop process and the change in the camera intrinsics.

Formally each flow prediction happens in a region of interest (RoI) within the original image, with size and location $B = [x, y, w, h]$. In our system the per-object flow map is predicted at a fixed size $w_r \times h_r$ using a RCNN based architecture as detailed in Section 4. For our view synthesis based loss functions we need to transform the image intrinsics $K = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$ into RoI specific versions.

The change only affects the intrinsic camera parameters and hence we need to compute a new intrinsic matrix K^j for each RoI j . The transformation ends up to be a displacement of the principal point and scaling of the focal length -

$$K^j = \begin{bmatrix} f_x w_r / w & 0 & (c_x - x) w_r / w \\ 0 & f_y h_r / h & (c_y - y) h_r / h \\ 0 & 0 & 1 \end{bmatrix}.$$

Note that we do not need bounding box associations between different viewpoints or time instants. We only compute detections for frame I_t^l and use a slightly expanded area as our RoI in frames that we warp to our reference frame for computing consistency losses in Eq. 1 and 3.

3.4. RoI Assembly for Full Frame Scene Flow

We assemble a complete scene flow from the object specific maps F^j . However, overlapping RoIs and certain RoIs may even contain multiple moving objects. Therefore we predict an object mask M^j for each RoI j in addition to F^j . The full 3D scene flow map F is computed as:

$$F = \sum_j M^j \odot F^j \quad (4)$$

We then use the full image flow map F with Eq. 1 and Eq. 3 for full image photometric and geometric losses. Note that the assembly procedure is fully differentiable and we are able to train instance masks $\mathcal{M} = \{M^1, \dots, M^j\}$ without any explicit mask supervision. We later use these instance masks (with flow) to identify moving objects (cf. Figure 6).

3.5. Full Learning Objective

We first state our full image synthesis based loss and then explain further priors we impose in our training loss. Our image synthesis loss function is based on four images I_t^l , I_t^r , I_{t+1}^l and I_{t+1}^r and can be split into three parts

$$\mathcal{L}^{\text{tot}} = \mathcal{L}^{\text{lr}} + \mathcal{L}^{\text{RoI}} + \mathcal{L}^{\text{t}} \quad (5)$$

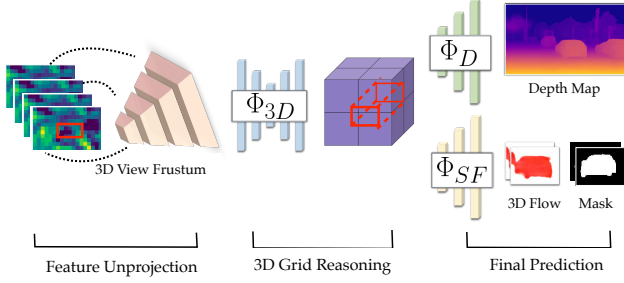


Figure 5. Network architecture. Our system predicts depth and instance-level 3D scene flow in a single forward pass. With extracted image features, we unproject features into a discretized view frustum grid, and then use a 3D CNN Φ_{3D} and finally perform prediction using depth Φ_D and scene flow Φ_{SF} decoders.

Where \mathcal{L}^{lr} is the loss for left-right consistency, \mathcal{L}^{RoI} is the RoI based loss function and \mathcal{L}^t is the full image based loss function on flow and depth over time. To state how the three parts are defined we introduce the notation $s \rightarrow t$ to indicate the warping from source s to target t .

$$\begin{aligned}\mathcal{L}^{lr} &= \mathcal{L}_{\text{photo}}(I_t^l, \hat{I}_t^{r \rightarrow l}) + \mathcal{L}_{\text{photo}}(I_{t+1}^l, \hat{I}_{t+1}^{r \rightarrow l}) \\ \mathcal{L}^{\text{RoI}} &= \sum_j \mathcal{L}_{\text{photo}}(I_t^{l,j}, \hat{I}_{t+1 \rightarrow t}^{l,j}) + \mathcal{L}_{\text{geo}}(D_t^{l,j}, \hat{D}_{t+1 \rightarrow t}^{l,j}, F_t^{l,j}) \\ \mathcal{L}^t &= \mathcal{L}_{\text{photo}}(I_t^l, \hat{I}_{t+1 \rightarrow t}^l) + \mathcal{L}_{\text{geo}}(D_t^l, \hat{D}_{t+1 \rightarrow t}^l, F_t^l)\end{aligned}\quad (6)$$

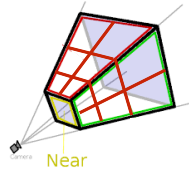
Beside the loss detailed above, we use additional priors such as smoothness for depth and flow while respecting discontinuities at boundaries [12]. Optionally, we use the classical stereo system ELAS [9] to compute an incomplete disparity map and use it for weak supervision with an L_1 loss.

4. Network Architecture

Figure 5 illustrates our network for scene flow, mask and depth prediction. We first talk about the 3D grid representation used to integrate the information from all images and then describe each component of the network.

4.1. 3D Grid Representation

In order to enable the network to reason about the scene geometry in 3D, we unproject the 2D features into a 3D grid [20]. A common discretization is to split a 3D cuboid volume of interest into equally sized voxels. This representation is used for 3D object shape reconstruction [40, 20]. However, it is not suitable for outdoor scenes with a large depth range, where we want to be more certain about foreground objects' geometry and motion, and allow increasing uncertainty with increasing depth in the 3D world. This lends to using the well known frustum shaped grid called matching cost volume or plane



sweep volume in classical (multi-view) stereo. In learning based stereo it has recently been used in [48]. The grid is discretized in image space plus an additional inverse depth ("nearness") coordinate, as shown in above image.

4.2. Network Components

Image Encoder. In the first stage the images are processed using a 2D CNN Φ_I , which outputs for each image a 2D feature map with c feature channels. The weights for this CNN are shared for all input frames - typically stereo frames at two time instants $\{I_t^l, I_t^r\}$ and $\{I_{t+1}^l, I_{t+1}^r\}$.

Unprojection. Using the 3D grid defined in Section 4.1, we lift the 2D information into the 3D space. We use the two left camera images as reference images $\{I_t^l, I_{t+1}^l\}$ and generate these 3D grids in both their camera coordinates. Each grid is populated with image features from all 4 images by projecting the grid cell centers into the respective images using the corresponding projection matrices [20]. We use the left images as reference frames as we predict disparity maps and scene flow from I_t^l to I_{t+1}^l .

Grid Pooling. The grids from the previous stage contain image features from all 4 frames. In order to combine the information from multiple frames we use two strategies. We use element-wise max pooling for features from left and right pairs and concatenate the features for different time instants in each grid cell. The motivation is that for stereo frames, there is no object motion and hence the feature should align well after unprojection. Thus a simple strategy of max pooling works well. Whereas for frames at different time instants, we expect motion in the scene and thus there would be misalignment where objects move. The output from this stage are two grids G_t^l and G_{t+1}^l .

3D Grid Reasoning. The next module Φ_{3D} processes the above two grids independently and generates output grids of the same resolution \tilde{G}_t^l and \tilde{G}_{t+1}^l . This module is implemented as a 3D encoder-decoder CNN module with skip connections following the U-Net architecture [35].

Output Modules. The final output is based on two CNN modules - one producing full frame depth for each reference image and one producing scene flow for each RoI in frame I_t . For each image I_i^l , with $i \in \{t, t+1\}$ we first collapse \tilde{G}_i^l (a 4D tensor) into a 3D tensor C_i^l by concatenating features in the depth dimension. As the grid is aligned with the reference image's camera, this corresponds to accumulating the features from various disparity planes at every pixel into a single feature. This tensor is further processed using ϕ_D to produce the full frame disparity map. The 3D flow prediction follows an RCNN [11] based architecture where given RoIs, we crop out corresponding regions C_t^l using an RoI align layer [15] and pass them to ϕ_{SF} which predict the scene flow and instance mask for each RoI. We also use skip connections from the image encoder in ϕ_D and ϕ_{SF} to

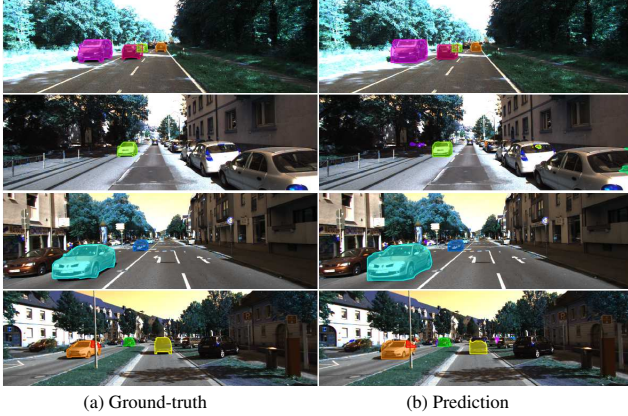


Figure 6. Qualitative results on our instance-level moving object mask prediction. Instances are color-encoded.

produce sharper predictions. The full frame scene flow map is created from the RoIs by pasting back as described in Section 3.4. The final outputs from our system are disparity maps D_t^l and D_{t+1}^l and a forward scene flow map F_t^l .

5. Experiments

We evaluate our instance-level 3d object motion and mask prediction on the KITTI 2015 scene flow dataset [29]. This is the only available dataset that contains real images together with ground-truth scene flow annotations. Following existing work [28, 49, 52, 12], we adopt the official 200 training images as test set. The official testing set is adopted for the final finetuning process. This is possible as we do not require the ground truth for training. All the related images in the 28 scenes covered by test data are excluded for training. Figure 6 and Figure 7 show some qualitative results.

Training details Our system is implemented using TensorFlow [1]. All models are optimized end-to-end using Adam [22] with a learning rate of 1×10^{-4} , decay rate of 0.5 and decay steps of 100000. During training, we randomly crop the input images in the horizontal direction to obtain patches with the size of 384×640 as input to the network. We set the output size of each RoI as 128×128 , we set the number of channels in the 3D grid to 64. The batch size is set as 1 to deal with flexible RoI number for training patch. For the image encoder, we finetune the first 4 convolutional layers from Inception ResNet V2 [37] pretrained on ImageNet. The rest of network is trained from scratch. We first train the depth prediction for 80K iterations on the KITTI raw dataset and then jointly train the depth and scene flow prediction for another 100k iterations. We finetune the model on the official testing set for another 120k iterations and use official 200 training images for comparison with other methods. The whole training process takes about 30 hours using a single NVIDIA Titan-X GPU.

5.1. Moving Object Speed and Direction Evaluation

Our method predicts 3D scene flow for each independently moving object. For each test image pair, ground-truth annotation of the disparity image at time t , the disparity image at time $t + 1$ warped into the first image’s coordinate frame and the 2D optical flow from time t to time $t + 1$ are provided. Using these GT annotations together with the estimated camera egomotion obtained from Libviso2 [10], we compute the 3D scene flow in the format of (x, y, z) for each image. To provide an instance-level analysis, we use the bbox detections [24], and find the dominant 3d flow for each object. As a result, we represent the motion direction and speed for each instance using a single 3d flow vector in the ground truth and all algorithms. We evaluate with the following metrics: the mean average error of the euclidean length of the 3d flow (speed), the mean average error of the angle of the 3d flow (motion direction) from the moving object pixels. For robustness to outliers we report the percentage of the mean average error below different thresholds. For comparison with other self-supervised flow and depth learning methods we need to reconstruct scene flow from depth and optical flow prediction. Geonet provides depthmaps with unknown scale factor and unflow does not estimate depth, we therefore use the depth results from Godard *et al.* [12]. As shown in Table 1, the average instance-level motion direction error of our method is less than 23° , about 15% smaller than the result obtained from the best self-supervised optical flow combined with the best self-supervised depth algorithm. In our prediction, about 75% of moving instances have an angular error below 15° .

5.2. Moving Object Instance Mask Evaluation

Our method can produce instance-level moving object segmentation from object bounding boxes and stereo videos. This is achieved without any instance mask ground truth supervision. We evaluate our predictions on the KITTI scene flow 2015 training split. The dataset provides an “Object map” which contains the foreground moving cars in each image. We use this motion mask as ground truth in our segmentation evaluation. Figure 6 shows some qualitative result of our moving object mask prediction. As shown in Table 2, we evaluate our mask prediction using the Intersection Over Union (IoU) metric. Specifically, We compute the mean image-level IoU which considers both moving object and static background and the mean instance-level IoU for only moving objects. Our method achieves highest IoU for mask prediction. As a baseline comparison, we use mask generated from SSD [24] 2D bounding box detections. Those masks contain both moving and static cars, thus it can only achieve an mean IoU of 0.34 for the full image mask. Even with the GT object movement information, it does not have tight object boundary and thus can only achieve a mean IoU of 0.655. This illustrates how

Method	AMAD \downarrow	AMAE \downarrow	AE $\leq 15^\circ\uparrow$	AE $\leq 30^\circ\uparrow$	SMAD \downarrow	SMAE \downarrow	SE $\leq 0.15\uparrow$	SE $\leq 0.3\uparrow$
GeoNet [49] + Godard [12]	6.98 $^\circ$	28.82 $^\circ$	62.93	77.16	0.256	0.503	0.351	0.554
UnflowC [28] + Godard [12]	5.96 $^\circ$	26.94 $^\circ$	64.87	77.58	0.240	0.471	36.21	58.62
Ours (no RoI consistency loss)	6.03 $^\circ$	29.34 $^\circ$	67.59	75.94	0.207	0.358	37.46	58.93
Our 3D scene flow	5.19$^\circ$	22.92$^\circ$	74.78	78.87	0.193	0.334	40.95	62.72

Table 1. Comparison of instance-level object motion in terms of motion direction(A) and speed (S). MAE denotes the mean average error, MAD denotes the median absolute deviation. The lower the better. We also report the percentage of the angle/speed error below different thresholds, where AE denotes the absolute angular error, SE denotes the absolute speed error. The higher the better.

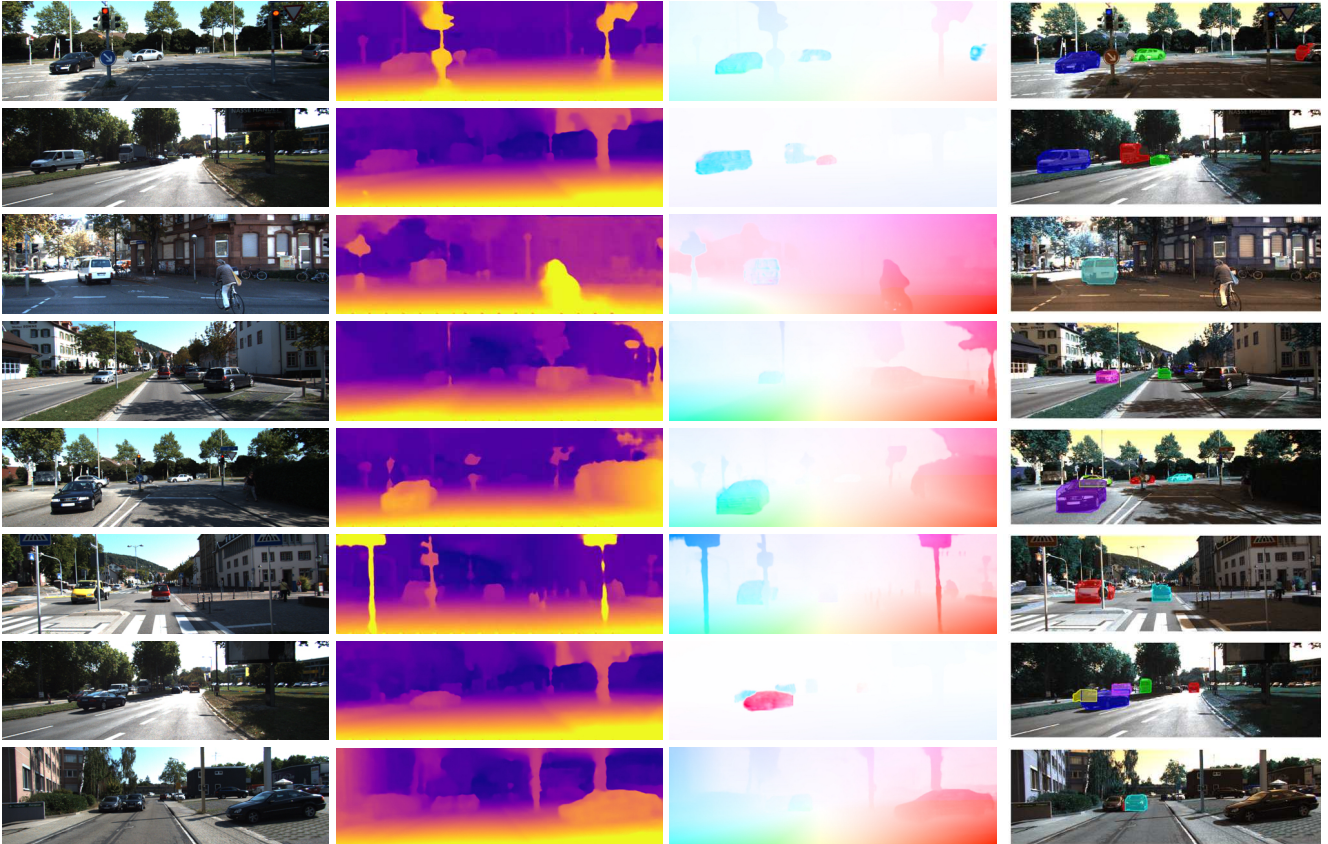


Figure 7. Qualitative results of our method. From left to right, reference image, depth, optical flow and instance-level moving object mask.

Method	Image IoU	Instance IoU
Zhou et al. [51]	0.380	-
Bounding box detections [24]	0.365	0.655
Our mask prediction	0.624	0.842

Table 2. Moving object mask evaluation. We report IoU number in both the full image and the moving instance bounding box.

our method effectively learns to determine which object is moving and identify an accurate instance segmentation for moving cars. We improve the result on both image-level and instance-level IoU. We also compare with Zhou *et al.* [51] which generates the foreground mask for all moving objects and occlusion region in the image. Their methods do not

provide instance-level information, hence we cannot obtain the instance-level IoU numbers.

5.3. Optical Flow Evaluation

An additional evaluation is to project our 3D flow predictions back to 2D to obtain the optical flow. As shown in Table 5, our method achieves the lowest EPE in both non-occluded regions and overall regions compared to other self-supervised methods. As a baseline comparison, we train a model without RoI consistency loss, which shows a decrease in performance. Optionally, we add an optical flow refinement sub-network, to further improve our optical flow result. The subnetwork is a unet which takes the warped image and the raw optical flow, together with orig-

Method	D1			D2			FL			ALL		
	bg	fg	bg+fg	bg	fg	bg+fg	bg	fg	bg+fg	bg	fg	bg+fg
EPC [47]	23.62	27.38	26.81	18.75	70.89	60.97	25.34	28.00	25.74			
EPC++ [26] (mono)	30.67	34.38	32.73	18.36	84.64	65.63	17.57	27.30	19.78	>30.67	>84.64	>65.63
EPC++ [26] (stereo)	22.76	26.63	23.84	16.37	70.39	60.32	17.58	26.89	19.64	>22.76	>70.39	>60.32
Godard <i>et al.</i> [12]	9.43	18.74	10.86	-	-	-	-	-	-	-	-	-
GeoNet [49]	-	-	-	-	-	-	43.54	48.24	44.26	-	-	-
Godard [12] + GeoNet flow	9.43	18.74	10.86	9.10	25.95	25.42	43.54	48.24	44.26	48.22	55.75	49.38
Ours	6.27	15.95	7.76	8.46	23.60	10.92	14.36	51.25	20.16	16.58	53.20	22.64

Table 3. Results on KITTI 2015 scene flow training split. All number shows the percentage of correctly predicted pixels. D1 denotes the disparity image at time t , D2 denotes the disparity image at time $t + 1$ warped into the first frame, FL denotes the 2D optical flow between the two time instances, fg denotes the foreground, and bg denotes the background.

Method	Binocular	Abs Rel	Sq Rel	RMSE
Godard <i>et al.</i> [12]	no	0.124	1.388	6.125
LIBELAS [9]	yes	0.186	2.192	6.307
Godard <i>et al.</i> [12]	yes	0.068	0.835	4.392
Ours	yes	0.065	0.699	3.896

Table 4. Results on the KITTI 2015 stereo training set of 200 disparity images. All learning-based methods are trained on KITTI raw dataset excluding the testing image sequences. The top half shows method which uses monocular image as input, the bottom half shows methods which use binocular images as input.

Method	Dataset	Non-occluded	All Regions
EpicFlow [34]	-	4.45	9.57
FlowNetS [5]	C+ST	8.12	14.19
FlowNet2 [17]	C+T	4.93	10.06
GeoNet [49]	K	8.05	10.81
DF-Net [52]	K+SY	-	8.98
UnFlowC [28]	K+SY	-	8.80
Ranjan <i>et al.</i> [33]	K	-	7.76
Ours	K	4.97	5.39
Ours (refined)	K	4.19	5.13

Table 5. Results on KITTI 2015 flow training set over non-occluded regions and overall regions. We use the average end-point error (EPE) metric to do the comparison. The classical method EpicFlow takes 16s per frame at runtime; The FlowNetS and FlowNet2 are learned with GT flow supervision. SY denotes SYNTHIA dataset [36], ST denotes Sintel dataset, C denotes FlyingChairs dataset, T denotes FlyingThings3D dataset. Numbers from other methods are directly taken from the paper.

inal image frames as input. This enables the network to further improve the optical flow prediction in a similar way as the architecture proposed in [32].

5.4. Depth Evaluation

To evaluate our depth prediction we use the KITTI 2015 stereo training set of 200 disparity images as test data and compare to other self-supervised learning and classical al-

gorithms in Table. 4. We compare to algorithms that take binocular stereo as input at test time. Our method achieves a higher accuracy as we input two consecutive binocular frames and our network also manages to match over time.

5.5. Scene Flow Evaluation

We compare other unsupervised method in the sceneflow subset by directly using their released results or running their released code. For this benchmark, a pixel is considered to be correctly estimated if the disparity or flow end-point error is ≤ 3 pixels or $\leq 5\%$. For scene flow this criterion needs to be fulfilled for two disparity maps and the flow map. As shown in Table 3, our method has an overall better accuracy than earlier self-supervised methods. Compared to classical approaches which optimize at test time our accuracy is still lower. However, test time optimization is in general prohibitively slow for real-time systems.

6. Conclusion

We presented a system to predict depth and object scene flow. Our network is trained using raw stereo sequences with off-the-shelf object detectors using image consistency as key learning objective. Our formulation is general and can be applied in any setting where a dynamic scene is imaged by multiple cameras - e.g. a multi-view capture system [13]. In future work, we would like to extend our system to integrate longer range temporal information. An emergent notion of objects to remove the dependence on pretrained object detectors is a further research direction. We also intend to explore general scenarios such as casual video captures using dual camera consumer devices and leverage large scale training for a truly general purpose depth and scene flow prediction system.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow:

- A system for large-scale machine learning. In *OSDI*, 2016. 6
- [2] Christoph Bregler, Aaron Hertzmann, and Henning Biermann. Recovering non-rigid 3d shape from image streams. In *CVPR*, 2000. 2
- [3] João Paulo Costeira and Takeo Kanade. A multibody factorization method for independently moving objects. *IJCV*, 1998. 2
- [4] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *NIPS*, 2014. 1, 2
- [5] P. Fischer, A. Dosovitskiy, E. Ilg, P. Häusser, C. Hazirbas, V. Golkov, P. v.d. Smagt, D. Cremers, and T. Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, 2015. 1, 2, 8
- [6] Katerina Fragkiadaki, Marta Salas, Pablo Arbelaez, and Jitendra Malik. Grouping-based low-rank trajectory completion and 3d reconstruction. In *NIPS*, 2014. 2
- [7] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *TPAMI*, 2010. 2
- [8] Ravi Garg, Vijay Kumar BG, Gustavo Carneiro, and Ian Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *ECCV*, 2016. 1, 2
- [9] Andreas Geiger, Martin Roser, and Raquel Urtasun. Efficient large-scale stereo matching. In *ACCV*, 2010. 5, 8
- [10] Andreas Geiger, Julius Ziegler, and Christoph Stiller. Stereoscan: Dense 3d reconstruction in real-time. In *Intelligent Vehicles Symposium (IV)*, 2011. 2, 6
- [11] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 5
- [12] C Godard, O Mac Aodha, and GJ Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, 2017. 1, 2, 5, 6, 7, 8
- [13] Lei Tan Lin Gui Bart Nabbe Iain Matthews Takeo Kanade Shohei Nobuhara Hanbyul Joo, Hao Liu and Yaser Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *ICCV*, 2015. 8
- [14] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. 2
- [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 5
- [16] Frédéric Huguet and Frédéric Devernay. A variational method for scene flow estimation from stereo sequences. In *ICCV*, 2007. 2
- [17] 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 *CVPR*, 2017. 1, 8
- [18] Mariano Jaimez, Mohamed Souiai, Javier Gonzalez-Jimenez, and Daniel Cremers. A primal-dual framework for real-time dense rgb-d scene flow. In *ICRA*, 2015. 2
- [19] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. Surfacenet: An end-to-end 3d neural network for multiview stereopsis. In *ICCV*, 2017. 2
- [20] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. In *NIPS*, 2017. 2, 5
- [21] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. In *ICCV*, 2017. 1, 2
- [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [23] Lubor Ladicky, Jianbo Shi, and Marc Pollefeys. Pulling things out of perspective. In *CVPR*, 2014. 2
- [24] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016. 4, 6, 7
- [25] H Christopher Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 1981. 2
- [26] 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. *arXiv preprint arXiv:1810.06125*, 2018. 2, 8
- [27] Nikolaus Mayer, Eddy Ilg, Philip Häusser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *CVPR*, 2016. 1, 2
- [28] Simon Meister, Junhwa Hur, and Stefan Roth. Unflow: Unsupervised learning of optical flow with a bidirectional census loss. *AAAI*, 2018. 6, 7, 8
- [29] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *CVPR*, 2015. 2, 6
- [30] Kemal Egemen Ozden, Kurt Cornelis, Luc Van Eycken, and Luc Van Gool. Reconstructing 3d trajectories of independently moving objects using generic constraints. *CVIU*, 2004. 2
- [31] Marc Pollefeys, Luc Van Gool, Maarten Vergauwen, Frank Verbiest, Kurt Cornelis, Jan Tops, and Reinhard Koch. Visual modeling with a hand-held camera. *IJCV*, 2004. 2
- [32] Anurag Ranjan and Michael J. Black. Optical flow estimation using a spatial pyramid network. In *CVPR*, 2017. 8
- [33] Anurag Ranjan, Varun Jampani, Kihwan Kim, Deqing Sun, Jonas Wulff, and Michael J Black. Adversarial collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation. *arXiv preprint arXiv:1805.09806*, 2018. 2, 8
- [34] Jerome Revaud, Philippe Weinzaepfel, Zaid Harchaoui, and Cordelia Schmid. Epicflow: Edge-preserving interpolation of correspondences for optical flow. In *CVPR*, 2015. 8
- [35] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 5
- [36] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M. Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *CVPR*, 2016. 8
- [37] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the

- impact of residual connections on learning. In *AAAI*, volume 4, page 12, 2017. 6
- [38] Tatsunori Tanai, Sudipta N Sinha, and Yoichi Sato. Fast multi-frame stereo scene flow with motion segmentation. In *CVPR*, 2017. 2
 - [39] Shubham Tulsiani, Tinghui Zhou, Alexei A Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *CVPR*, 2017. 2
 - [40] Shubham Tulsiani, Tinghui Zhou, Alexei A. Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *CVPR*, 2017. 5
 - [41] Shimon Ullman. Maximizing rigidity: The incremental recovery of 3-d structure from rigid and nonrigid motion. *Perception*, 1984. 2
 - [42] Sundar Vedula, Simon Baker, Peter Rander, Robert Collins, and Takeo Kanade. Three-dimensional scene flow. In *ICCV*. IEEE, 1999. 2
 - [43] Sudheendra Vijayanarasimhan, Susanna Ricco, Cordelia Schmid, Rahul Sukthankar, and Katerina Fragkiadaki. Sfmnet: Learning of structure and motion from video. Technical report, arXiv:1704.07804, 2017. 2
 - [44] Christoph Vogel, Konrad Schindler, and Stefan Roth. Piecewise rigid scene flow. In *ICCV*, 2013. 2
 - [45] Zhou Wang, Alan C Bovik, Hamid R Sheikh, Eero P Simoncelli, et al. Image quality assessment: from error visibility to structural similarity. *TIP*, 2004. 3
 - [46] Andreas Wedel, Clemens Rabe, Tobi Vaudrey, Thomas Brox, Uwe Franke, and Daniel Cremers. Efficient dense scene flow from sparse or dense stereo data. In *ECCV*, 2008. 2
 - [47] Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, and Ram Nevatia. Every pixel counts: Unsupervised geometry learning with holistic 3d motion understanding. *arXiv preprint arXiv:1806.10556*, 2018. 2, 8
 - [48] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *ECCV*, 2018. 2, 5
 - [49] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *CVPR*, 2018. 1, 2, 4, 6, 7, 8
 - [50] Chang Yuan and Gerard Medioni. 3d reconstruction of background and objects moving on ground plane viewed from a moving camera. In *CVPR*, 2006. 2
 - [51] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, 2017. 2, 7
 - [52] Yuliang Zou, Zelun Luo, and Jia-Bin Huang. Df-net: Unsupervised joint learning of depth and flow using cross-task consistency. In *ECCV*, 2018. 2, 4, 6, 8