

Depth from Videos in the Wild: Unsupervised Monocular Depth Learning from Unknown Cameras

Ariel Gordon^{1,2}, Hanhan Li², Rico Jonschkowski^{1,2}, Anelia Angelova^{1,2}

¹Robotics at Google, ²Google AI

{gariel, uniqueness, rjon, anelia}@google.com

Abstract

We present a novel method for simultaneously learning depth, egomotion, object motion, and camera intrinsics from monocular videos, using only consistency across neighboring video frames as a supervision signal. Similarly to prior work, our method learns by applying differentiable warping to frames and comparing the result to adjacent ones, but it provides several improvements: We address occlusions geometrically and differentiably, directly using the depth maps as predicted during training. We introduce randomized layer normalization, a novel regularizer, and we account for object motion relative to the scene. To the best of our knowledge, our work is the first to learn the camera intrinsic parameters, including lens distortion, from video in an unsupervised manner, thereby allowing us to extract accurate depth and motion from arbitrary videos of unknown origin at scale. We evaluate our results on the Cityscapes, KITTI, and EuRoC MAV datasets, establishing new state of the art on depth prediction and odometry, and demonstrate qualitatively that depth prediction can be learned from a collection of YouTube videos. The code is publicly available¹.

1. Introduction

Estimating 3D structure and camera motion from video is a key problem in computer vision. Traditional approaches to this problem rely on identifying the same points in the scene in multiple consecutive frames, then solving for a 3D structure and camera motion that is maximally consistent across those frames [23]. But such correspondences between frames can only be established for a subset of all pixels, which leaves the problem of estimating depth underdetermined. As commonly done with inverse problems, the gaps are filled based on assumptions of continuity, planarity, etc.

¹github.com/google-research/google-research/tree/master/depth_from_video_in_the_wild

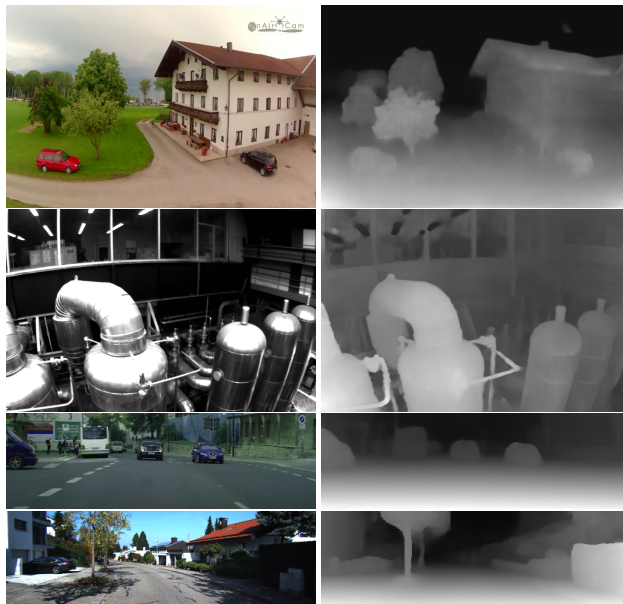


Figure 1. Qualitative results of our approach for learning depth from videos of unknown sources, which is enabled by simultaneously learning the camera extrinsic and intrinsic parameters. Since our method does not require knowing the camera parameters, it can be applied to any set of videos. All depth maps (visualized on the right, as disparity) were learned from raw videos without using any camera intrinsics groundtruth. From top to bottom: frames from YouTube8M [1], from EuRoC MAV dataset [5], from Cityscapes [7] and from KITTI [11].

Rather than specifying these assumptions manually, deep learning is able to obtain them from data. Wherever information is insufficient to resolve ambiguities, deep networks can produce depth maps and flow fields by generalizing from prior examples they have seen. Unsupervised approaches allow learning from raw videos alone, using similar consistency losses as traditional methods but optimizing them during training. At inference, the trained networks are able to predict depth from a single image and egomotion from pairs or longer sequences of images.

As research in this direction gained traction [47, 10, 12, 33, 24, 34], it became clear that object motion is a major obstacle because it violates the assumption of a static scene. Several directions have been proposed to address the issue [44, 40], including leveraging semantic understanding of the scene through instance segmentation [6]. Occlusions have been another limiting factor, and lastly, in all prior work in this direction, the intrinsic parameters of the camera had to be provided.

This work addresses the problems above and, as a result, reduces supervision and improves the quality of depth and motion prediction from unlabeled videos. First, we show that a deep network can be trained to predict the *intrinsic* parameters of the camera, including lens distortion, in an *unsupervised* manner from the video itself (see Fig. 1). Second, we are the first in this context to address occlusions directly, in a geometric way, from the predicted depth as it is. Lastly, we substantially reduce the amount of semantic understanding needed to address moving elements in the scene: Instead of segmenting every instance of a moving object and tracking it across frames [6], we need a single mask that covers pixels that *could* belong to a moving object. This mask can be as rough as a union of rectangular bounding boxes. Obtaining such a rough mask is a much simpler problem than instance segmentation and can be solved more reliably with existing models.

In addition to these qualitative advances, we conduct an extensive quantitative evaluation of our method and find that it establishes a new state of the art on multiple widely used benchmark datasets. Pooling datasets together, a capability which is greatly advanced by our method, proves to enhance quality. Finally, we demonstrate for the first time that depth and camera intrinsics prediction can be learned on YouTube videos, which were captured with multiple different cameras, each with unknown and generally different intrinsics.

2. Related work

Estimating scene depth is an important task for robot navigation and manipulation. Historically much research has been devoted to it, including a large bodies of research on stereo, multi-view geometry, and active sensing [29, 21, 9]. Recently, learning-based approaches for dense depth prediction have gained attention [9, 22, 19, 45]. In these, scene depth is predicted from input RGB images and the depth estimation function is learned using supervision provided by a sensor, such as a LiDAR. Similar approaches are used for other dense predictions such as surface normals [8, 38].

Unsupervised depth learning. Unsupervised learning of depth, where the only supervision is obtained from the monocular video itself and no depth sensors are needed, has also been popularized recently [47, 10, 12, 33, 24, 34, 44].

Garg et al. [10] introduced joint learning of depth and egomotion. Zhou et al. [47] demonstrated a fully differentiable approach where depth and egomotion are predicted jointly by deep neural networks. Techniques were developed for the monocular [33, 42, 24, 41, 44, 35, 6] and binocular [12, 33, 40, 46, 35, 46] settings. In the latter, it was shown that monocular depth quality at inference is improved when stereo inputs are used during training. Other methods learn directly the stereo disparity [17, 18, 43]. Other novel techniques include the use of motion [41, 35, 44, 6, 40].

Learning from images or videos from unknown cameras. This is an active research field, focusing on single or multi-view images [2, 30, 20]. It is especially hard for internet photos due to the diversity of input sources and lack of the camera parameters, as shown by Li et al. [20]. Our work makes a step in addressing this challenge by learning camera intrinsics for videos in the wild.

Occlusion aware learning. Multiple approaches disconnected from geometry have been proposed for handling occlusions in the context of optical flow [36, 15, 25]. Differentiable mesh rendering [26, 16] adopts a geometric approach to occlusions. In the context of learning to predict depth and egomotion, occlusions were addressed via a learned explainability mask [47], by penalizing the minimum reprojection loss between the previous frame or the next frame into the middle one, and through optical flow [40]. In the latter context, we are the first to address occlusions in a direct geometric approach via a differentiable loss.

Learning of intrinsics. Learning to predict the camera intrinsics has mostly been limited to strongly supervised approaches. The sources of groundtruth vary: Workman et al. [37] use focal lengths estimated employing classical 1D structure from motion. Yan et al. [39] obtain the focal length based on EXIF. Bogdan et al. [4] synthesize images from panoramas using virtual cameras with known intrinsics, including distortion. To our knowledge, our approach is the only one that learns the camera intrinsics in an unsupervised manner directly from video, jointly with depth, egomotion and object motion.

3. Preliminaries

Similarly to prior work [47, 12, 44, 32], the backbone of our method is the equation that ties together two adjacent video frames using a depth map (z) and the camera matrix K . Eq. 1 describes the shift in a pixel position p due to a rotation matrix R and a translation vector t :

$$z'p' = KRK^{-1}zp + Kt \quad (1)$$

p' and z' are the new homogeneous coordinates of the pixel and the new depth.

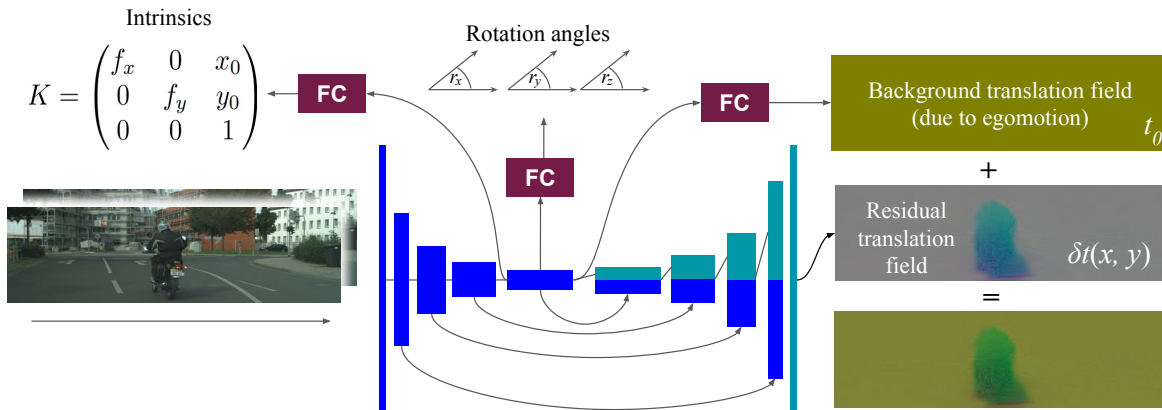


Figure 2. A schematic of our motion-prediction network. The network receives two images as input. A stack of convolutions creates a bottleneck, from which fully-connected (FC) network heads predict the intrinsics, rotation angles and translation vector components of the background (due to egomotion). A series of decoder layers predict a residual translation field, which predicts motion of objects with respect to the scene (color coding is explained in Fig. 4). A mask (not shown in the figure, see Eq. 2) multiplies by zero the residual translation field at all pixels that don't belong to possibly-moving objects. The residual translation field is then added to the background translation to obtain a total translation field. The images are taken from Cityscapes. A separate network (now shown) predicts depth from a single image.

Using z , R , and t as predicted by deep networks, Eq. 1 is used to warp one video frame onto the other. The result is then compared to the actual other frame, and the differences constitute the main component of the training loss. The premise is that through penalizing the differences, the networks will learn to correctly predict z , R , and t .

4. Method

In this work we propose simultaneous learning of depth, egomotion, object motion, and camera intrinsics from monocular videos. A motion-prediction network predicts camera motion, motion of every pixel with respect to the background, and the camera intrinsics: focal lengths, offsets and distortion. A second network predicts depth maps. By imposing consistency across neighboring frames as a loss, the networks simultaneously learn to predict depth maps, motion fields and the camera intrinsics. To apply this loss only in unoccluded pixels, we estimate occlusions geometrically based on estimated depth maps. We regularize motion fields based on masks that indicate which pixels might belong to moving objects, obtained from a pretrained segmentation or object detection network.

4.1. Networks and losses

Networks Depth is predicted from a *single image* by a UNet [28] encoder-decoder network with a ResNet 18 base and a softplus activation ($z(\ell) = \log(1 + e^\ell)$) to convert the logits (ℓ) to depth (z).

A second network (shown in Fig. 2) predicts camera motion, a dense residual translation representing motion of objects relative to the scene, as well as the camera intrinsics,

from *two consecutive images*. Further details about the network are given in the Supplementary Material (SM).

Losses Based on the estimated depth map, camera intrinsics, rotation, and the translation field, we warp the first frame to match the second one and compare those using two losses: 1) a structural similarity (SSIM) loss and 2) the sum of $L1$ distances for the color channels, following Casser et al. [6]. Additionally, we impose a cycle consistency loss on the motion field by estimating both forward and backward motion, which we obtain by applying the networks on the frames in normal and in reversed order. Since those motion estimates at corresponding pixels should be opposite, we define an L2 loss on the relative deviation from opposite rotation and translation. Additionally, we apply spatial L1 smoothness losses for the depth and motion fields, a temporal L1 smoothness loss for depth, and an L2 weight regularization term. More details are given in the SM.

4.2. Occlusion-aware consistency

When the camera and / or objects move, areas in the scene that were visible in one frame may become occluded in another, and vice versa. Photometric consistency cannot be enforced in pixels that correspond to these areas. Given a depth map and a motion field in one frame, one could actually detect where occlusion is about to occur, and exclude the occluded areas from the consistency loss. Detecting occluded pixels requires some sort of reasoning about the connectivity of the surface represented by the depth map, and z-buffering. Keeping the mechanism differentiable and efficient enough for a training loop, may pose a challenge.

We therefore take a different approach, as illustrated in Fig. 3. For each pixel (i, j) in the source frame, the pre-

dicted depth z_{ij} and the camera intrinsic matrix are used to obtain the respective point in space, (x_{ij}, y_{ij}, z_{ij}) . The point is moved in space according to the predicted motion field. In particular, the depth changes to z' . The new spatial location is reprojected back onto the camera frame, and falls at some generally-different location (i', j') on the target frame. i' and j' are generally non-integer. Therefore obtaining the depth on the target frame at (i', j') , $z_{i', j'}^t$, requires interpolation.

Occlusions happen at (i', j') where z' becomes multivalued. At such points, color and depth consistency should be applied only to the visible branch of z' , that is, the branch where z' is smaller. If the source and target frames are nearly consistent, the visible branch will be close to target depth at (i', j') , $z_{i', j'}^t$. The way we propose to pick the visible branch is to include in the losses only points (i', j') where $z_{i', j'}^s \leq z_{i', j'}^t$. In other words, only if a transformed pixel on the source frame lands in front of the depth map in the target frame, do we include that pixel in the loss. This scheme is not symmetrical with respect to interchanging the source and target frames, which is why we always apply it in a symmetrized way: We transform the source onto the target, calculate the losses, and then switch the roles of source and target. Fig. 3 illustrates the method.

The losses described in Sec. 4.1 are invoked in an “occlusion-aware” manner, as described in this section, except for SSIM. For the latter, we handle occlusions by replacing all averaging operations by a weighted averaging, where the weight of a pixel is a decreasing function of the depth error in that pixel. The exact expression is given in the SM.

4.3. Regularization

Semantic regularization of the translation field Eq. 1 can propagate frame inconsistency losses to z , R and t at every pixel. However without further regularization, they remain greatly underdetermined. While continuity of z , R and t is a powerful regularizer, we found that further regularization helps significantly. In particular, we impose constancy of R throughout the image, and allow t to deviate from a constant value only at pixels that are designated as *possibly mobile*. Unlike in prior work [6], instance segmentation and tracking are not required, as all we need is a single “possibly mobile” mask $m(x, y)$. We write

$$t(x, y) = t_0 + m(x, y)\delta t(x, y), \quad (2)$$

where t_0 and $\delta t(x, y)$ are the background motion (due to camera motion) and residual motion (due to object motion).

We show in ablation experiments that $m(x, y)$ can be as rough as a union of bounding boxes (see Fig. 4). In addition, an $L1$ smoothing operator is applied to $t(x, y)$.

Randomized layer normalization In our experiments, we

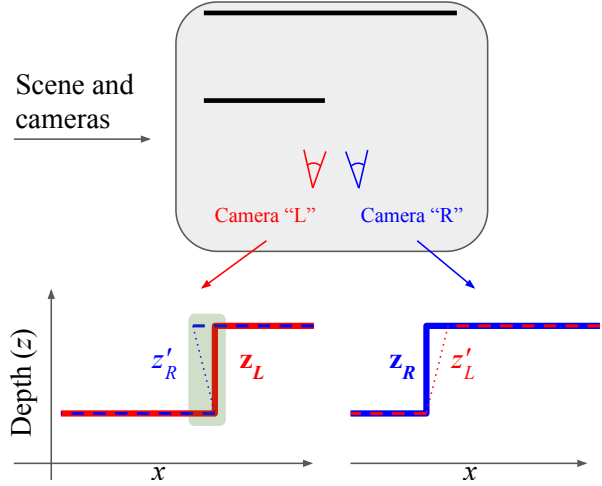


Figure 3. An illustration of our proposed method for handling occlusions. At the top we show a two-dimensional “scene”, consisting of two straight surfaces, one partially occluding the other. Two cameras, left (“L”) and right (“R”), are observing the scene. Our method is *monocular*, so these represent two locations of the same camera that moved, and “left” and “right” are used for convenience. At the bottom, the depth map observed by each camera is illustrated as a solid line on the respective side (z_L and z_R). A dashed line shows the depth map obtained from warping one view onto the other (z'_R and z'_L). The warped depth map can become a multivalued function, which indicates occlusions (see the green-shaded rectangle). To handle that, we apply the photometric and geometric losses only at pixels where $z'_R \leq z_L$ and $z'_L \leq z_R$. When the depth maps and motion estimation are correct, the loss in this scheme would indeed evaluate to zero.

observed the following anomalous behavior in relation to Batch Normalization (BN):

- Evaluation metrics were consistently better when running inference at the “training mode” of BN. That is, instead of long-term averaged means and variances, the ones obtained from the image itself during inference were used², rendering batch normalization more similar to layer normalization [3].
- As we increased the batch size at training, the eval accuracy was consistently worse and worse, no matter how we scaled the learning rate.

These two anomalies led to the conclusion that BN is acting like Layer Normalization (LN) [3], and for each item in the batch, the others act as a source of noise. Replacing BN by LN with multiplicative Gaussian noise led to improved evaluation metrics and allowed increasing the batch size (accompanied with a linear increase of the learning rate [14]) without loss and even with a slight improvement in the evaluation metrics.

²Even at batch size 1, there would still be means and variances over the spatial dimensions.



Figure 4. Use of a “possibly mobile” mask to regularize the translation field. An object detection network identifies all instances of objects that are capable of motion, such as pedestrians, cyclists and cars. The union of the bounding boxes comprises the “possibly mobile” mask, within which the translation field is allowed to vary. The top picture, from Cityscapes, illustrates the mask, and the bottom one is the translation field predicted by the network (x, y, z coded as RGB). The **golden background** corresponds to motion in the negative z direction, as the entire scene is moving towards the camera. The **greenish silhouette** is the cyclist moving slightly to the left and slightly towards the camera. Note that the network carves the silhouette out of a rough mask.

4.4. Learning the intrinsics

Through Eq. 1, photometric consistency losses across neighboring frames provide a supervision signal for K , but only when there is a nonzero camera rotation between them $R \neq \mathbb{1}$. Indeed, when $R = \mathbb{1}$ Eq. 1 is reduced to $z'p' = zp + Kt$, which means that the loss depends on K only through the product Kt . Kt , however, can be perfectly correct even if K and t are incorrect. In fact, for any (non-singular) \tilde{K} , there exists a $\tilde{t} = \tilde{K}^{-1}Kt$, such that $\tilde{K}\tilde{t} = Kt$. Especially when K and t are predicted by two heads stemming from the same network, the latter can “escape” the supervision signal as long as it predicts a K and a t whose product is correct.

Fortunately, rotations can provide a supervision signal for K . Eq. 3 (derived in the SM) ties the tolerance with which the focal lengths can be determined (δf_x and δf_y , denominated in pixels) to the amount of camera rotation that occurs between the two:

$$\delta f_x < \frac{2f_x^2}{wsr_y}; \quad \delta f_y < \frac{2f_y^2}{hsr_x}. \quad (3)$$

r_y and r_x are the rotation angles along the respective axes in radians, w and h are the image width and height respectively, and $s = \max(h, w)$.

5. Experiments

In this section, we evaluate our method on depth prediction, odometry estimation, and the recovery of camera

intrinsics across a range of diverse datasets.

5.1. Datasets

KITTI The KITTI dataset is collected in urban environments and is the main benchmark for depth and egomotion estimation. It is accompanied with a LIDAR sensor, which is used for evaluation only. We use standard splits into training, validation, and test sets, commonly referred to as the Eigen split. 39835 training examples are used from KITTI.

Cityscapes The Cityscapes dataset is a more recent urban driving dataset, which we use for both training and evaluation. It is a more challenging dataset with many dynamic scenes. With a few exceptions [27, 6] it has not been used for depth estimation evaluation. It has 38675 training examples. We use depth from the disparity data for evaluation on a standard evaluation set of 1250 samples [27, 6].

EuRoC Micro Aerial Vehicle Dataset The EuRoC Micro Aerial Vehicle (MAV) Dataset [5] is a very challenging dataset collected by an aerial vehicle indoors. While the data contains a comprehensive suite of sensor measurements, including stereo pairs, IMU, accurate Leica laser tracker ground truth, Vicon scene 3d scans, and camera calibration, we only used the monocular videos for training. Since the camera has significant lens distortion, this is an opportunity to test our method for learning lens distortions.

YouTube8M videos To demonstrate that depth can be learned on videos in the wild from unknown cameras, we collected videos from the YouTube8M dataset [1]. From the 3079 videos in YouTube8M that have the label “quadcopter”, human raters selected videos that contain significant amount of footage from a quadcopter. Naturally, the videos were taken with different unknown cameras, with varying fields of view and varying degrees of lens distortion. The YouTube8M IDs are listed in the SM.

5.2. Depth

Since monocular methods can only estimate depth up to a global scale factor, we follow the field’s common practice [47] of normalizing out the scale factor based on the medians of predicted and groundtruth depths (code).

KITTI Table 1 summarizes the evaluation results on the KITTI Eigen partition of a model trained on KITTI. The metrics are the ones defined in Zhou et al. [47]. Only the best methods and the first three metrics are displayed in Table 1, the rest are given in the SM. As seen in the table, we improve on the state-of-the-art results. More importantly, we observed that learned intrinsics, rather than given ones, consistently help performance.

Cityscapes Table 2 summarizes the evaluation metrics of models trained and tested on Cityscapes. We follow the established protocol by previous work, using the disparity for

Method	M	Abs Rel	Sq Rel	RMSE
Zhou [47]		0.208	1.768	6.856
Yang [42]		0.182	1.481	6.501
Mahjourian [24]		0.163	1.240	6.220
LEGO [41]	✓	0.162	1.352	6.276
GeoNet [44]	✓	0.155	1.296	5.857
DDVO [35]		0.151	1.257	5.583
Godard [13]		0.133	1.158	5.370
Struct2Depth [6]	✓	0.141	1.026	5.291
Yang [40]		0.137	1.326	6.232
Yang [40]	✓	0.131	1.254	6.117
Ours:				
Given intrinsics	✓	0.129	0.982	5.23
Learned intrinsics	✓	0.128	0.959	5.23

Table 1. Evaluation of depth estimation of our method, with given and learned camera intrinsics, for models trained and evaluated on KITTI, compared to other monocular methods. The depth cutoff is always 80m. The “M” column is checked for all models where object motion is taken into account.

evaluation [6, 27]. Since this is a very challenging benchmark with many dynamic objects, very few approaches have evaluated on it. As seen in Table 2, our approach outperforms previous ones and benefits from learned intrinsics.

Method	M	Abs Rel	Sq Rel	RMSE
Pilzer [27]		0.440	5.713	5.443
Struct2Depth [6]	✓	0.145	1.736	7.279
Ours:				
Given intrinsics	✓	0.129	1.35	6.96
Learned intrinsics	✓	0.127	1.33	6.96

Table 2. Evaluation of depth estimation of models trained on Cityscapes on the cityscapes test set, with a depth cutoff of 80m, and comparison to prior art.

Cityscapes + KITTI Being able to learn depth without the need for intrinsics opens up the opportunity for pooling videos from any data source. Figure 5 shows the results of pooling Cityscapes and KITTI datasets and evaluating on either one. In this experiment the intrinsics are assumed unknown and are learned. Training jointly on both datasets improves the depth metrics even beyond the best results obtained on either dataset separately. This is a key result which demonstrates the impact of our method to leverage data sources of potentially unlimited size.

Cityscapes + KITTI: ablation experiments Table 3 summarizes the results of ablation experiments we ran in order to study the impact of each of the techniques described in this paper on the end results. In order to reduce the number of combinations of results, in all experiments the training set was Cityscapes and KITTI mixed together. Each model was evaluated on both Cityscapes and KITTI separately.

Using a union of bounding boxes as a “possibly-mobile” mask (as depicted in Fig. 4) is found to be as good as us-

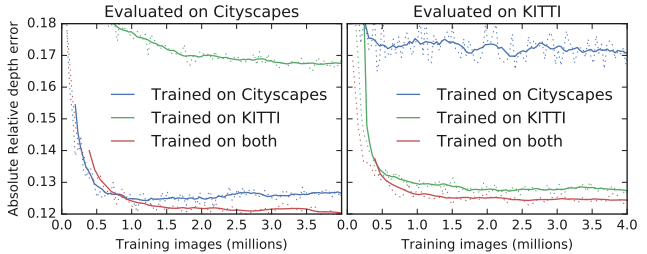


Figure 5. Depth prediction on Kitti and Cityscapes when training on each dataset and on both. Joining the two datasets improves the results on both. Learning intrinsics allows to similarly pool many datasets, even if they are from unknown cameras.

Experiment	Abs Rel depth	
	CS	KITTI
Our algorithm	0.121	0.124
Boxes instead of segmentation masks	0.120	0.125
w/o occlusion-aware loss	0.127	0.126
w/o object motion	0.172	0.130
w/o randomized layer normalization	0.124	0.127

Table 3. Ablation experiments on depth estimation. In all experiments the training set was Cityscapes (CS) and KITTI combined, and we tested the model on Cityscapes (CS) and KITTI (Eigen partition) separately. Each row represents an experiment where one change was made compared to the main method, as described in the “Experiment” row. Smaller numbers are better.

ing segmentation masks, which makes our technique more broadly applicable. Object motion estimation is shown to play a crucial role, especially on Cityscapes, which is characterized by more complex scenes with more pedestrians and cars. Randomized LN is shown to be superior to standard BN, and lastly – occlusion-awareness improves the quality of depth estimation, especially on Cityscapes, which has richer scenes with more occlusions. Figure 6 visualizes the type of artifacts that occlusion-aware losses reduce.

EuRoC MAV Dataset We further use the EuRoC MAV Dataset to evaluate depth. We selected also a very challenging out-of-sample evaluation protocol in which we trained on the “Machine room” sequences and tested on the “Vicon Room 2 01”, which has 3D groundtruth. Table 4 reports the results. The SM details how depth ground truth was generated from the provided Vicon 3D scans.

Abs R	Sq R	RMS	lgRMS	a_1	a_2	a_3
0.332	0.389	0.971	0.396	0.420	0.743	0.913

Table 4. Evaluation of depth estimation for EuRoC dataset, no prior results are available for this dataset. $a_i = \delta < 1.25^i$.

5.3. YouTube Videos

To demonstrate that depth can be learned on collections of videos in the wild for which the camera parameters are

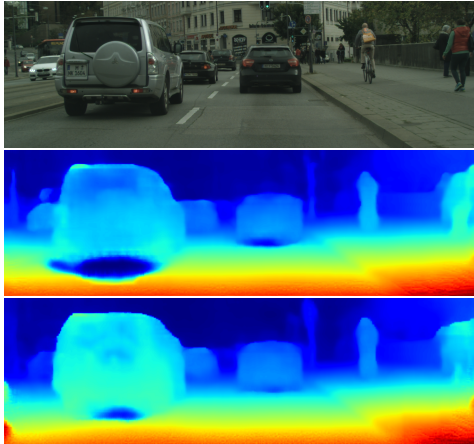


Figure 6. Illustration of the effect of occlusion aware losses. The center and bottom images are inferred disparity maps obtained from the image at the top. In the center image, the model was trained without occlusion-aware losses. At areas that become disoccluded, under cars, occlusion aware loss is shown to reduce artifacts. The top image belongs to the Cityscapes test set and is one of images whose depth prediction metrics were hurt the most upon removing occlusion aware losses.

not known, and differ across videos, we trained our model on the YouTube8M videos described in Sec. 5.1. Figure 7 visualizes the results. We note that this dataset is very challenging as it features objects at large ranges of depth.



Figure 7. Frames from from YouTube and learned disparity maps. The camera intrinsics are learned.

5.4. Camera intrinsics evaluation

In evaluating our method for learning camera intrinsics, two separate questions can be asked. First, how good is the supervision signal that cross-frame consistency provides for the camera intrinsics. Second is how good a deep network is in learning them and generalizing to test data.

Quality of the supervision signal. To evaluate the supervision signal for the intrinsics, we represented each intrinsic parameter as a separate learned variable. This is suitable for the EuRoC dataset, since it was captured with the same camera throughout. Each of the 11 subsets of EuRoC were trained in a separate experiment, until the intrinsics converged, yielding 11 independent results. The resulting

sample mean and standard deviation for each intrinsic parameter are summarized in Table 5. All parameters agree with groundtruth within a few pixels. Since the groundtruth values were not accompanied by tolerances, it is hard to tell whether or not the differences are within tolerance. The learned disparity maps are shown in Fig. 8.

Quantity	Learned	GT
Horizontal focal length (f_x)	253.7 ± 1.1	250.2
Vertical focal length (f_y)	265.4 ± 1.3	261.3
Horizontal center (x_0)	189.0 ± 0.9	187.2
Vertical center (y_0)	132.2 ± 1.1	132.8
Quadratic radial distortion	-0.267 ± 0.003	-0.283
Quartic radial distortion	0.064 ± 0.002	0.074

Table 5. Camera intrinsics learned on the EuRoC datasets. Learning of depth, egomotion and intrinsics was done separately on each of the 11 datasets, using monocular images (“cam0”) only. Constancy of the intrinsics throughout each dataset separately was imposed, and statistics (mean and standard deviation) for each intrinsic parameter were collected across the results. The groundtruth (GT) was adjusted to an image size of (256×384).



Figure 8. Frames from EuRoC [5] and the corresponding learned disparity maps. The camera intrinsics are learned.

Learning and generalization Prior work [37, 39, 4] has shown that deep networks can learn and generalize camera intrinsics in a strongly supervised setting. In our setting, the camera intrinsics and motion are predicted by the same network and are thus correlated. In other words, the losses imposed on the motion / intrinsics network only impose correctness of the intrinsics within the limits of Eq. 3.

We evaluated our model’s predictions of the intrinsics on the KITTI odometry series. The model was trained on the union of the Cityscapes and KITTI training sets, which significantly differ in their typical focal lengths. Figure 9 shows the scatter plot of the predicted f_x as function of the predicted r_y . The predictions fall within the limits imposed by Eq. 3. While Table 5 indicates a high-quality supervision signal for the intrinsics, Fig. 9 shows that when the intrinsics and motion are predicted by the same network, the latter learns them “in aggregate”, only to the extent needed for the depth-prediction network to learn depth.

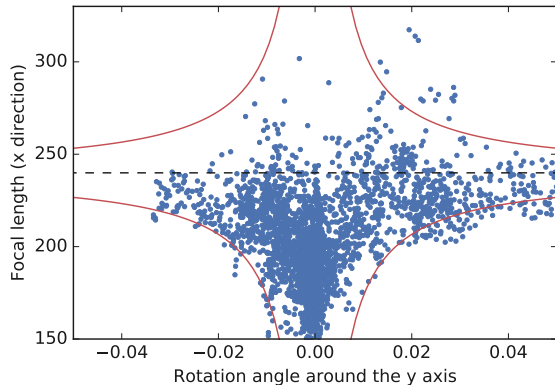


Figure 9. Predicted f_x as function of the predicted r_y for all images in the KITTI odometry sequences 09 and 10. The dashed line is groundtruth; red curves show tolerance limits imposed by Eq. 3.

5.5. Odometry

We evaluated our egomotion prediction on the KITTI sequences 09 and 10. The common 5-point Absolute Trajectory Error (ATE) metric [47, 6, 44, 13] measures local agreement between the estimated trajectories and the respective groundtruth. However assessing the usefulness for a method for localization requires evaluating its accuracy in predicting location. A common metric for localization is average relative translational drift t_{rel} [31, 46] – the distance between the predicted location and the groundtruth location divided by distance traveled and averaged over the trajectory. Table 6 summarizes both metrics, demonstrating the improvements our method achieves on both.

When evaluating for odometry, the most naive way is to calculate the inference of egomotion for every pair of adjacent frame. That leads to the red “learned intrinsics” curve in Fig. 10. However it is also possible to make an inference-time correction if we know the intrinsics of the camera at inference time. In that case, one can leverage the fact that for small errors in the rotation angles and focal lengths, $r_x f_y$ and $r_y f_x$ are approximately constant (Eq. SM8). Therefore if the network predicted r'_y and f'_x for a given pair of images, and we know the true focal length f_x , we can correct our estimate of r_y to $r'_y f'_x / f_x$. This is the procedure invoked in generating the “Learned and corrected intrinsics” curve in Fig. 10, and the respective entry in Table 6. The trajectories and metrics obtained with learned intrinsics with inference time correction and with given intrinsics are similar. Both notably improve prior art, which is especially prominent for localization, as the t_{rel} metric indicates.

6. Conclusions

This work addresses important challenges for unsupervised learning of depth and visual odometry through geo-

Metric	Seq. 09		Seq. 10	
	ATE	t_{rel}	ATE	t_{rel}
Zhou [47]	0.021	17.84%	0.020	37.91%
GeoNet [44]	0.012	/	0.012	/
Zhan [46]	/	11.92%	/	12.45%
Mahjourian [24]	0.013	/	0.012	/
Struct2depth [6]	0.011	10.2%	0.011	28.9%
Ours, with intrinsics:				
Learned	0.012	7.5%	0.010	13.2%
Learned & corrected	0.010	2.7%	0.007	6.8%
Given	0.009	3.1%	0.008	5.4%

Table 6. Absolute Trajectory Error (ATE) [47] and average relative translational drift (t_{rel}) [31] on the 09 and 10 KITTI odometry sequences. Our method with both learned and given intrinsics is compared to prior work.

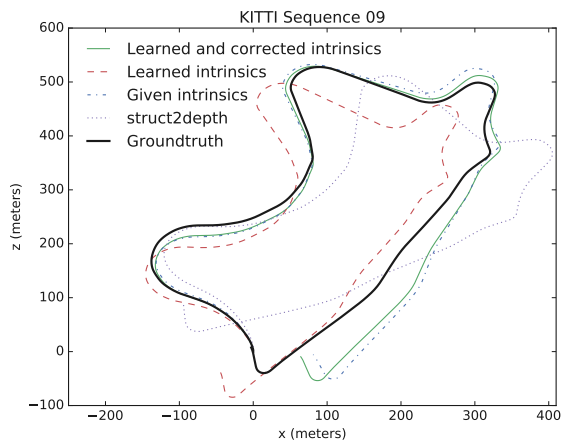


Figure 10. Predicted location on the KITTI odometry sequence 09, generated by models trained on KITTI, with given intrinsics and with learned intrinsics (with and without inference time correction), compared to groundtruth and to struct2depth results.

metric handling of occlusions, a simple way of accounting for object motion, and a novel form of regularization. Most importantly, it takes a major step towards leveraging the vast amounts of existing unlabeled videos for learning depth estimation: Through unsupervised learning of the camera intrinsic parameters, including lens distortion, it enables, for the first time, learning depth from raw videos captured by unknown cameras.

7. Acknowledgements

We Thank Roland Siegwart for the permission to use EuRoC dataset, Sergey Ioffe for enlightening discussions about batch normalization, Jason Su and Kwea Koi for spotting an error in the derivation in an earlier version of the manuscript, Kurt Konolige for consultation and critical reading of the manuscript, and Chad Richards for copyediting the paper.

References

- [1] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudeendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *CoRR*, abs/1609.08675, 2016. 1, 5
- [2] Sameer Agarwal, Noah Snavely, Ian Simon, Steven M. Seitz, and Richard Szeliski. Building Rome in a day. *ICCV*, 2009. 2
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 4
- [4] Oleksandr Bogdan, Viktor Eckstein, Francois Rameau, and Jean-Charles Bazin. Deepcalib: A deep learning approach for automatic intrinsic calibration of wide field-of-view cameras. In *Proceedings of the 15th ACM SIGGRAPH European Conference on Visual Media Production*, CVMP '18, pages 6:1–6:10, New York, NY, USA, 2018. ACM. 2, 7
- [5] Michael Burri, Janosch Nikolic, Pascal Gohl, Thomas Schneider, Joern Rehder, Sammy Omari, Markus W Achtelik, and Roland Siegwart. The euroc micro aerial vehicle datasets. *The International Journal of Robotics Research*, 2016. 1, 5, 7
- [6] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Unsupervised learning of depth and ego-motion: A structured approach. In *AAAI-19*, 2019. 2, 3, 4, 5, 6, 8, 14, 16
- [7] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 1
- [8] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. *ICCV*, 2015. 2
- [9] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *NIPS*, 2014. 2
- [10] Ravi Garg, Gustavo Carneiro, and Ian Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. *ECCV*, 2016. 2
- [11] 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. 1
- [12] Clément Godard, Oisín Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. *CVPR*, 2017. 2
- [13] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel Brostow. Digging into self-supervised monocular depth estimation. *arxiv.org/pdf/1806.01260*, 2018. 6, 8, 14, 16
- [14] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large mini-batch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 4
- [15] Joel Janai, Fatma Guney, Anurag Ranjan, Michael Black, and Andreas Geiger. Unsupervised learning of multi-frame optical flow with occlusions. *ECCV*, 2018. 2
- [16] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [17] 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 *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 2
- [18] Sameh Khamis, Sean Fanello, Christoph Rhemann, Adarsh Kowdle, Julien Valentin, and Shahram Izadi. Stereonet: Guided hierarchical refinement for real-time edge-aware depth prediction. In *The European Conference on Computer Vision (ECCV)*, September 2018. 2
- [19] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. *arXiv:1606.00373*, 2016. 2
- [20] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. *CVPR*, 2018. 2
- [21] Fayao Liu, Chunhua Shen, and Guosheng Lin. Deep convolutional neural fields for depth estimation from a single image. *CVPR*, 2015. 2
- [22] F. Liu, C. Shen, G. Lin, and I. Reid. Learning depth from single monocular images using deep convolutional neural fields. *PAMI*, 2015. 2
- [23] David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, 60(2):91–110, Nov. 2004. 1
- [24] Reza Mahjourian, Martin Wicke, and Anelia Angelova. Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5667–5675, 2018. 2, 6, 8, 14, 16
- [25] Michal Neoral, Jan Sochman, and Jir Matas. Continual occlusions and optical flow estimation. *ECCV*, 2018. 2
- [26] Thu H Nguyen-Phuoc, Chuan Li, Stephen Balaban, and Yongliang Yang. Rendernet: A deep convolutional network for differentiable rendering from 3d shapes. In *Advances in Neural Information Processing Systems*, pages 7902–7912, 2018. 2
- [27] Andrea Pilzer, Dan Xu, Mihai Marian Puscas, Elisa Ricci, and Nicu Sebe. Unsupervised adversarial depth estimation using cycled generative networks. *3DV*, 2018. 5, 6, 14
- [28] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells, and Alejandro F. Frangi, editors, *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pages 234–241, Cham, 2015. Springer International Publishing. 3
- [29] Ashutosh Saxena, Jamie Schulte, and Andrew Y. Ng. Depth estimation using monocular and stereo cues. In *Proceedings*

- of the 20th International Joint Conference on Artificial Intelligence, IJCAI'07, pages 2197–2203, San Francisco, CA, USA, 2007. Morgan Kaufmann Publishers Inc. 2
- [30] Johannes L. Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. *CVPR*, 2016. 2
- [31] Jared Shamwell, Sarah Leung, and William Nothwang. Vision-aided absolute trajectory estimation using an unsupervised deep network with online error correction. 10 2018. 8, 16
- [32] Shubham Tulsiani, Tinghui Zhou, Alexei A. Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. *CVPR*, 2017. 2
- [33] Benjamin Ummerhofer, Huizhong Zhou, Jonas Uhrig, Nikolaus Mayer, Eddy Ilg, Alexey Dosovitskiy, and Thomas Brox. Demon: Depth and motion network for learning monocular stereo. *CVPR*, 2017. 2
- [34] Sudheendra Vijayanarasimhan, Susanna Ricco, Cordelia Schmid, Rahul Sukthankar, and Katerina Fragkiadaki. Sfm-net: Learning of structure and motion from video. *arXiv:1704.07804*, 2017. 2
- [35] Chaoyang Wang, Jose Miguel Buenaposada, Rui Zhu, and Simon Lucey. Learning depth from monocular videos using direct methods. *CVPR*, 2018. 2, 6, 14
- [36] Yang Wang, Yi Yang, Zhenheng Yang, Liang Zhao, Peng Wang, and Wei Xu. Occlusion aware unsupervised learning of optical flow. *CVPR*, 2018. 2
- [37] Scott Workman, Connor Greenwell, Menghua Zhai, Ryan Baltenberger, and Nathan Jacobs. DeepFocal: A Method for Direct Focal Length Estimation. In *International Conference on Image Processing*, 2015. 2, 7
- [38] Abhinav Gupta Xiaolong Wang, David F. Fouhey. Designing deep networks for surface normal estimation. *CVPR*, 2015. 2
- [39] Han Yan, Yu Zhang, Shunli Zhang, Sicong Zhao, and Li Zhang. Focal length estimation guided with object distribution on focalens dataset. *Journal of Electronic Imaging*, 26(3):1 – 14 – 14, 2017. 2, 7
- [40] Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, and Ram Nevatia. Every pixel counts: Unsupervised geometry learning with holistic 3d motion understanding. *arxiv.org/pdf/1806.10556*, 2018. 2, 6, 14
- [41] Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, and Ram Nevatia. Lego: Learning edge with geometry all at once by watching videos. *CVPR*, 2018. 2, 6, 14
- [42] Zhenheng Yang, Peng Wang, Wei Xu, Liang Zhao, and Ramakant Nevatia. Unsupervised learning of geometry with edge-aware depth-normal consistency. *arXiv:1711.03665*, 2017. 2, 6, 14
- [43] Yao Yao, Zixin Luo, Shiwei Li, and Tian Fang and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. *ECCV*, 2018. 2
- [44] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. *CVPR*, 2018. 2, 6, 8, 14, 16
- [45] Zifeng Wu Yuanzhouhan Cao and Chunhua Shen. Estimating depth from monocular images as classification using deep fully convolutional residual networks. *CoRR:1605.02305*, 2016. 2
- [46] 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. *CVPR*, 2018. 2, 8, 16
- [47] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G. Lowe. Unsupervised learning of depth and ego-motion from video. *CVPR*, 2017. 2, 5, 6, 8, 14, 16