

# Supplementary Material for Self-Supervised Scale Recovery for Monocular Depth and Egomotion Estimation

Brandon Wagstaff, and Jonathan Kelly

TABLE I: Network architecture of our pose encoder network. **k** refers to the kernel size, **s** refers to the stride, **chns** refers to the input/output channels.

Pose Network					
Layer	k	s	chns	input	Activation
conv1	3	1	8/16	img+flow	ReLU
conv2	3	2	16/32	conv1	ReLU
conv3	3	3	32/64	conv2	ReLU
conv4	3	2	64/128	conv3	ReLU
conv5	3	2	128/256	conv4	ReLU
conv6	3	2	256/512	conv5	ReLU
conv7	3	2	512/1024	conv6	ReLU
avgpool				conv7	-
transconv	1	1	1024/3	avgpool	-
rotconv	1	1	1024/3	avgpool	-

## I. NETWORK DETAILS

### A. Pose Encoder Network

Our pose network is shown in Table I. ReLU activations are applied after every convolutional layer (except for the final transconv and rotconv layers, as these must be able to regress negative values). The initial number of input channels is eight, as we use two RGB images, and the optical flow estimate between them.

### B. Depth Encoder Network

We use a standard Resnet18 [1] architecture to extract a feature representation from a single image.

### C. Depth Decoder Network

Our decoder architecture is similar in spirit to Monodepth2 [2] and DNet [3]: the main architecture is based on the decoder from Monodepth2, but we incorporate “dense connected prediction” (DCP) layers from DNet to combine features from different scales. We used three scales for our experiments.

At each layer of the baseline decoder, the feature space is upsampled by a factor of two within each upconv layer, and is added to a previous encoder output via a skip connection. See Table II for more details. The extracted features from the final three layers are passed into the DCP layers; see Table III for more information. Here, the number of channels from the final three iconv layers are reduced to eight using a convolutional layer, and are passed through disparity prediction layers, which reduce the number of channels to one via a 2D convolution of stride one; a sigmoid activation

TABLE II: Network architecture of our pose encoder network. Upsampling (represented with a  $\uparrow$ ) was done prior to applying the upconv convolution, and was done with nearest neighbour interpolation.

Depth Decoder					
Layer	k	s	chns	input	activation
upconv5	3	1	512/256	$\uparrow$ enc5	ELU
iconv5	3	1	256/256	upconv5	ELU
upconv4	3	1	256/128	$\uparrow$ iconv5	ELU
iconv4	3	1	128/128	upconv4 + enc4	ELU
upconv3	3	1	128/64	$\uparrow$ iconv4	ELU
iconv3	3	1	64/64	upconv3 + enc3	ELU
upconv2	3	1	64/64	$\uparrow$ iconv3	ELU
iconv2	3	1	64/64	upconv2 + enc2	ELU
upconv1	3	1	64/32	$\uparrow$ iconv2	ELU
iconv1	3	1	32/32	upconv1	ELU

TABLE III: Description of the DCP layers, which merge lower level features with higher level features prior to applying a disparity prediction.

DCP Layers					
Layer	k	s	chns	input	activation
dcp3	3	1	64/8	iconv3	ELU
disp3	3	1	8/1	dcp3	sigmoid
dcp2	3	1	64/8	iconv2	ELU
disp2	3	1	16/1	$\uparrow$ dcp3, dcp2	sigmoid
dcp1	3	1	32/8	iconv1	ELU
disp1	3	1	24/1	dcp3, $\uparrow$ dcp2, $\uparrow$ dcp1	sigmoid

is applied to constrain the network output to be within [0,1]. To produce the inverse depth estimate at each of these scales, the following transformation is applied:

$$\mathbf{D}_t^{-1} = \frac{1}{D_{min}} + \left( \frac{1}{D_{max}} - \frac{1}{D_{min}} \right) * \mathbf{out}, \quad (1)$$

which constrains the depth estimates to be in the range  $[D_{min}, D_{max}]$ .

### D. Plane Segmentation Network

Our plane segmentation network uses a pretrained Resnet18 Encoder. Following this, our decoder is described in Table IV

## II. TRAINING DETAILS

For all experiments, we use an image resolution of 192x640, and whiten all input images using the ImageNet [4] statistics. All images are augmented by randomly flipping images, and modifying the brightness, contrast, saturation, and hue within the ranges  $\pm 0.2, \pm 0.2, \pm 0.2, \pm 0.1$  respectively (the same modification is made to all the target and source images within a training sample). For training on KITTI, we use sequences 00, 02, 06–08, 11, 13–16, 19, and leave out 05 for validation, and 09–10 for testing.

TABLE IV: Plane Segmentation Decoder details.

Plane Decoder					
Layer	k	s	chns	Input	activation
upconv5	3	1	512/256	$\uparrow$ enc5	ELU
iconv5	3	1	256/256	upconv5	ELU
upconv4	3	1	256/128	$\uparrow$ iconv5	ELU
iconv4	3	1	128/128	upconv4 + enc4	ELU
upconv3	3	1	128/64	$\uparrow$ iconv4	ELU
iconv3	3	1	64/64	upconv3 + enc3	ELU
upconv2	3	1	64/64	$\uparrow$ iconv3	ELU
iconv2	3	1	64/64	upconv2 + enc2	ELU
upconv1	3	1	64/32	$\uparrow$ iconv2	ELU
iconv1	3	1	32/32	upconv1	ELU
plane1	3	1	32/1	iconv1	sigmoid

Currently, we separate the training of our plane segmentation network from the training of our depth and egomotion networks, because an accurate depth prediction is required to train the plane segmentation network (we use it to extract 3D coordinates from each image coordinate in the image). Therefore, we initially train an *unscaled* depth and egomotion network (with the our baseline, unscaled loss function), and use the unscaled depth predictions from this network to train the ground plane segmentation network. Our final model was chosen by selecting the one that resulted in the lowest validation loss. The hyperparameters for the plane segmentation network training were  $\lambda_{plane} = 25$ , and  $\lambda_{reg} = 0.05$ .

Following the training of the plane segmentation network, we train our *scaled* depth and egomotion networks with our scale recovery loss (along with the traditional photometric reconstruction losses described in the paper). We initialized our networks by training for one epoch with the unscaled (baseline) loss (eq. 7), since our scale recovery loss requires reasonable depth estimates to estimate the scale factor. After one epoch, we incorporated the scale recovery loss term (eq. 17) to begin resolving metric scale. Each training sample consists of three images: a target image, and two source images that are temporally adjacent to the target. We adopt several training improvements Monodepth2 [2] to improve the overall network training. Automasking is used to ignore pixel regions that do not change in intensity between frames, and the per-pixel minimum reprojection loss is computed to account for occluded and out-of-view pixels. Additionally, we use the same multi-scale training scheme (i.e., we evaluate our training loss at each of the depth network’s three depth resolutions, but upsample each scale to full resolution before computing the reconstruction loss). In addition to warping each source frame to the target frame, we do the inverse, by computing each source image depth, and warping the target frame to the source frames (and compute the photometric reconstruction loss for each). For our depth consistency loss, we warp each source depth to the target frame, and compare these to the target depth.

We set  $\lambda_P = 1$ ,  $\lambda_S = 0.05$ ,  $\lambda_{DC} = 0.14$ ,  $\lambda_{PC} = 5$ ,  $\lambda_{DS} = 0.02$ ,  $\lambda_{TS} = 0.6$ ,  $D_{min} = 1.8$ ,  $D_{max} = 60$ . The ground truth camera height was  $h_{gt} = 1.70$  meters.

### III. GROUND PLANE EXTRACTION

For Experiment 1 (scale factor evaluation), we used a RANSAC-based ground plane extraction method to estimate the per-frame scale factor. This technique was also used in Experiment 2 (visual odometry evaluation) to perform online rescaling (to reproduce DNet). Here, we describe the procedure that was used to estimate the camera height over the ground plane.

The ground plane was extracted with a RANSAC-based [5] plane fitting procedure. For each image, we segmented the lower-middle region that generally corresponded to a road plane (pixels within  $[\frac{4}{7}H, H]$ , and  $[\frac{1}{6}W, \frac{5}{6}W]$ , and determined the 3D coordinates for each pixel by reprojecting them using the camera intrinsics and the depth network prediction. During each RANSAC iteration, we fit a plane with three sampled points, and computed the number of inlier plane pixels (a pixel was considered an inlier if the distance to the plane was found to be less than 0.02 m). 350 RANSAC iterations were performed, and the ground plane that resulted in the highest number of inliers was chosen. To determine the camera height over the ground plane, we evaluated the offset  $d$  of each inlier point  $\mathbf{p}$  ( $\mathbf{p}^T \tilde{\mathbf{n}} = d$ ), and considered the median offset to be the estimated camera height.

### IV. RETRAINING EXPERIMENT

For Oxford Robotcar [6], [7] training, we kept all training conditions the same, with the exception of training for eight epochs only (due to there being a larger number of images in the training dataset), and adjusting the camera height to  $h_{gt} = 1.52$  meters.

The Robotcar models were trained on sequences 2014-11-18-13-20-12, 2015-07-08-13-37-17, 2015-07-10-10-01-59, 2015-08-12-15-04-18; we use the first and second subsequences of 2014-11-18-13-20-12 for validation and testing respectively.

### REFERENCES

- [1] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proc. IEEE Conf. Comput. Vision Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [2] C. Godard, O. Aodha, M. Firman, and G. Brostow, “Digging into self-supervised monocular depth estimation,” in *Proc. IEEE Conf. Comput. Vision Pattern Recognition (CVPR)*, 2019, pp. 3828–3838.
- [3] F. Xue, G. Zhuo, Z. Huang, W. Fu, Z. Wu, and M. A. Jr, “Toward hierarchical self-supervised monocular absolute depth estimation for autonomous driving applications,” *arXiv preprint arXiv:2004.05560*, 2020.
- [4] “ImageNet: A large-scale hierarchical image database,” in *Proc. IEEE Conf. Comput. Vision Pattern Recognition (CVPR)*.
- [5] M. Fischler and R. Bolles, “Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography,” *J. Commun. ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [6] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, “1 Year, 1000km: The Oxford RobotCar Dataset,” *Int. J. Robot. Res. (IJRR)*, vol. 36, no. 1, pp. 3–15, 2017. [Online]. Available: <http://dx.doi.org/10.1177/0278364916679498>
- [7] W. Maddern, G. Pascoe, M. Gadd, D. Barnes, B. Yeomans, and P. Newman, “Real-time kinematic ground truth for the oxford robotcar dataset,” *arXiv preprint arXiv: 2002.10152*, 2020. [Online]. Available: <https://arxiv.org/pdf/2002.10152>