

Monocular depth and motion estimation using neural networks

Erik Örjehag

Master of Science Thesis in Electrical Engineering

Monocular depth and motion estimation using neural networks

Erik Örjehag

LiTH-ISY-EX--YY/NNNN--SE

Supervisor: **Gustav Häger**
ISY, Linköpings universitet

Examiner: **Per-Erik Forssén**
ISY, Linköpings universitet

*Division of Computer Vision
Department of Electrical Engineering
Linköping University
SE-581 83 Linköping, Sweden*

Copyright © 2020 Erik Örjehag

Sammanfattning

Det här som vi har hållit på med är jätteviktigt faktiskt och det vi gjort blev bara sååå bra. Kanske inte helt otippat, men det glass är sååå gott!

Förresten har vi blivit bäst på att skriva rapporter, så nu ska ska vi inte gå in närmare på några detaljer såhär i sammanfattningen.

Abstract

If your thesis is written in English, the primary abstract would go here while the Swedish abstract would be optional.

Acknowledgments

Vi tycker alla har varit så himla goa hela den här långa och tuffa tiden i våra liv.

Linköping, Januari 2020
NN och MM

Contents

Notation	xi
1 Introduction	1
1.1 Contributions	1
1.2 Motivation	2
1.3 Research Questions	2
1.4 Delimitations	2
2 Background	3
2.1 Convolutional neural networks	3
3 Related work	5
3.1 Unsupervised Monocular Depth Estimation with Left-Right Consistency [4]	5
3.2 Unsupervised Learning of Depth and Ego-Motion from Video [13]	6
3.3 SuperDepth: Self-Supervised, Super-Resolved Monocular Depth Estimation [11]	7
3.4 3D Packing for Self-Supervised Monocular Depth Estimation [6]	7
3.5 Digging Into Self-Supervised Monocular Depth Estimation [5]	7
3.6 Self-Supervised 3D Keypoint Learning for Ego-motion Estimation [12]	8
4 Method	9
4.1 Unsupervised depth and ego motion learning	9
4.1.1 Sequence datasets	9
4.1.2 Architectures for depth and ego motion CNNs	10
4.1.3 Differentiable depth image warping	10
4.1.4 Loss functions	12
4.1.5 Handling occlusions	14
4.1.6 Handling model limitations	15
4.2 Unsupervised keypoint learning	16
4.2.1 Score encoder	18
4.3 Unsupervised consensus maximization	18

4.4 Evaluation	22
4.4.1 Depth	23
4.4.2 Ego motion	24
4.4.3 Keypoints	24
4.4.4 Consensus maximization	24
5 Results	25
6 Discussion	29
Bibliography	31

Notation

NÅGRA MÄNGDER

Notation	Betydelse
\mathbb{N}	Mängden av naturliga tal
\mathbb{R}	Mängden av reella tal
\mathbb{C}	Mängden av komplexa tal

FÖRKORTNINGAR

Förkortning	Betydelse
ARMA	Auto-regressive moving average
PID	Proportional, integral, differential (regulator)

1

Introduction

This thesis aims to show how new neural network techniques can be used to predict depth and relative camera motion for an image sequence captured by a monocular RGB camera. Imagine closing one eye and looking out into the world, it is trivial as a grownup human to detect motion and estimate how the head moves in relation to what is seen. Calculating camera movement from an image sequence is a well studied problem and is usually done by finding corresponding features in the images and calculating (using projective geometry) which camera movement can give rise to such correspondences and their relative movement between frames in the sequence.

Recent research has shown that it's possible to predict depth and relative motion from a sequence of images taken with a monocular RGB camera, up to a unknown scale factor. The training data is a sequence of unlabeled images with a small relative motion, for example looking out from the front window of a moving car. Given a target view and a new nearby view it's possible to train depth and pose predicting CNNs jointly using a combined loss function. The depth and pose predictions are used to warp nearby views to the target view and the loss is based on the similarity achieved after warping.

1.1 Contributions

This thesis does an analysis, and a performance comparison of different techniques described in current research papers. Training and testing is done on new datasets.

1.2 Motivation

Localization by only using visual input is highly desirable in robotics applications due to the low hardware cost and power consumption of using cameras compared to, for example, 3D lidars. Obtaining labeled data can be a tedious task, for that reason this thesis will focus on unsupervised learning on unlabeled data.

1.3 Research Questions

1. What ideas from previous work can be combined and what are the performance gains if any?
2. How well do previous methods work on new datasets not tested in the original papers? In what ways do the results break down?
3. What alteration to previous loss functions can be made to improve results?
4. How does the amount of training data affect the results? Can the performance of the methods from previous research be improved simply by training on more data?
5. Can the results from the original papers be replicated in the PyTorch framework?

1.4 Delimitations

The visual localization problem can be solved using, for example, a stereoscopic camera or a time of flight camera. But this thesis will only explore the use of a monoscopic, non depth sensing, RGB camera, because it enables applications where hardware cost is a big factor.

In a full SfM pipeline, both 3D-reconstruction, bundle adjustment and loop-closure detection are usually done as well. This will not be part of this thesis project.

2

Background

2.1 Convolutional neural networks

The central method used in this project is a deep learning algorithm called convolutional neural networks (CNN for short). A CNN architecture can successfully capture the spatial dependencies in an image through convolutional filtering operations with kernels of learnable weights and biases.

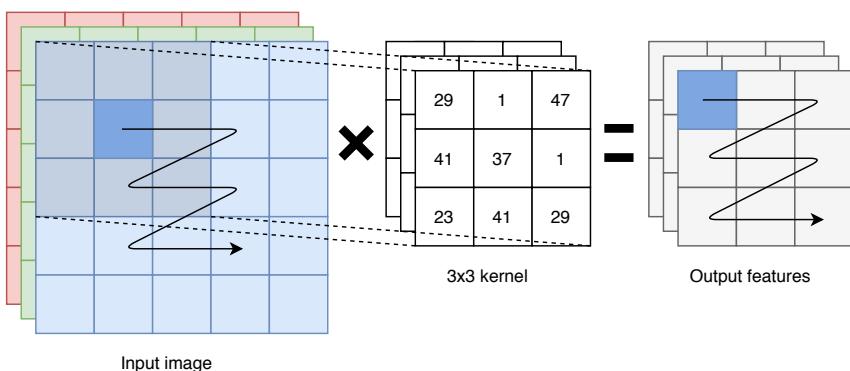


Figure 2.1: Convolutional filtering operation with a 3 channel RGB image and 3x3 kernel

In Figure 2.1 a convolutional filtering operation over an image is illustrated. The matrix kernel is moved in a row by row pattern and is multiplied by a patch of the image to get a value for the output cell.

In order to form a deep neural network multiple filtering operations are chained sequentially with nonlinear activation functions between them. A deep neural network usually consists of many such layers of filtering operations and activation functions. This concept is illustrated in Figure 2.2.

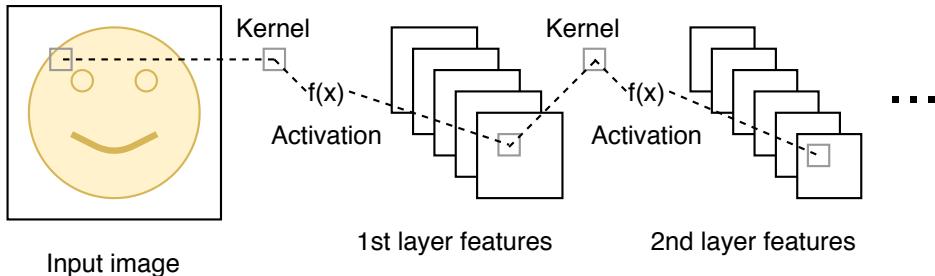


Figure 2.2: Multiple layers chained with each other to form a deep network

The activation functions need to be differentiable because the derivatives are used in the learning process. Some common nonlinear activation functions are shown in Figure 2.3.

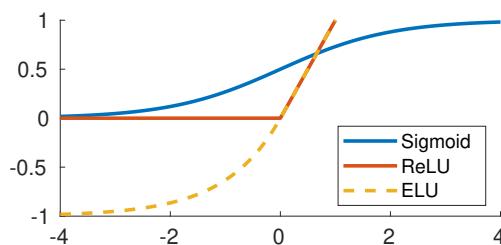


Figure 2.3: A few common nonlinear activation functions

The features in the deep neural networks are used to formulate a loss function. The loss defines an objective that we want the network to learn. This means that the process of learning becomes a task of updating the filter kernels in such a way that the loss function is minimized. If the loss function is decreasing during the training process it means that the network is learning. It is crucial that the loss function describes the problem accurately, otherwise the network will not learn the correct behavior. The weights in the network are updated using a method called back propagation. The algorithm computes the gradient of the loss function with respect to the weights in the network for a single input example from the training data using the derivative chain rule, which can be done very efficiently. Updating the weights to minimize the loss function can then be done using gradient descent. As the network is fed with more input examples from the training data, the network slowly learns the correct weights that minimizes the loss function as desired.

3

Related work

In this chapter some of the important contributions of previous research papers are summarized.

3.1 Unsupervised Monocular Depth Estimation with Left-Right Consistency [4]

In a paper from 12 Apr 2017 the authors present MonoDepth, with an implementation available in Tensorflow on github. In this work the depth is predicted using an encoder-decoder type network, but the relative motion between frames is not estimated at all. The KITTI dataset provides stereo image pairs which are used during training, and the relative transformation between the left and right cameras is known. Using only the left image as input to the network both disparity maps for the left and right images are predicted. The two disparity maps are used to project the left image into the view of the right image and vice versa. This can be seen as a precursor to the papers discussed later which uses only a monocular image sequence to train a depth predicting network and pose predicting network jointly.

In the loss function, the left image is reconstructed from the right image and vice versa. This is possible using the disparity maps and known baseline and focal length of the cameras. The L1 norm of the per pixel photometric error as well as SSIM are computed and added to the loss. An additional loss term encourages the left disparity to be equal to the right disparity projected into the left camera viewpoint. Because the photometric error does not work well on low textured regions an edge aware smoothness term is added to propagate the depth values from

nearby areas in the disparity map. The method produces metrically accurate results because the baseline and focal length of the cameras are known.

3.2 Unsupervised Learning of Depth and Ego-Motion from Video [13]

In a paper from CVPR 21 July 2017 the authors present SfMLearner with an official implementation in Tensorflow on github. Additional replications in the PyTorch and Chainer frameworks are also available.

Contrary to the paper in section 3.1 only a monocular sequence of images from the KITTI dataset is used during training. In the stereo case the relative pose between the left and right cameras was known, but with this monocular dataset the pose between subsequent frames is unknown. The authors train a pose predicting and depth predicting network simultaneously with a joint loss function to solve this problem.

During training 3 subsequent frames are considered at a time. The frame I_{t-1} and frame I_{t+1} are called the source frames and the frame I_t is called the target frame. The target frame is input to the depth network which estimates a disparity map. The two source frames are fed through a pose estimating network one after each other together with the target frame to find the relative transformations $T_{t \rightarrow t-1}$ and $T_{t \rightarrow t+1}$.

The authors add an explainability mask to the photometric error term to account for errors in the model. The view synthesis formulation implicitly assumes that the scene does not contain moving objects, that there are no occlusions between the target and source frames, and that the surfaces are Lambertian so that the photo-consistency error of RGB values is meaningful. In order to predict the explainability mask an additional CNN is used. The network has no explicit supervisory signal but is encouraged to be non-zero with an regularization term using a cross-entropy loss with a constant label 1 at each pixel location. This makes the network minimize the view synthesis objective but is allowed some slack due to factors not considered by the model. In later work it was shown that the explainability mask does not help to improve results that much and is often ignored.

To tackle the problem with textured areas and non Lambertian surfaces, a smoothness term is used. An edge aware smoothness term was not used like in previous work[4] but a penalty on the second order gradient of the depth map was used instead. This unfortunately makes edges very fuzzy in the results compared to other methods.

3.3 SuperDepth: Self-Supervised, Super-Resolved Monocular Depth Estimation [11]

In this paper published 3 Oct 2018 the authors train a depth estimating network on only stereo image pairs. After the depth predicting network has been trained the pose network from section 3.2 is trained, using results from the already trained depth predicting network. This means that the depth and pose networks are not trained with a joint loss function but are instead trained separately.

The main contributions from this paper is the proposal to use supixel-convolution. The method additionally uses differentiable flip-augmentation to remove edge artifacts on the left and right edges of the depth map seen in previous work using stereo image pairs during training.

To handle occluded pixels between the left and right images an occlusion regularization loss term is added to encourage background depths (low disparity).

3.4 3D Packing for Self-Supervised Monocular Depth Estimation [6]

The authors present PacknetSfM in their paper on 6 Dec 2019. The main contribution is a new network architecture with packing and unpacking blocks replacing down and up sampling. The new packing blocks use space to depth transformations and 3d convolutions. The authors claim that the new packing and unpacking blocks are better at preserving resolution than standard down and up-sampling. As the method is not adopted in later work, it is unclear if the new packing and unpacking approach is actually any good.

The second contribution is a loss on the camera velocity that makes it possible for the monocular depth estimation to be metrically accurate. The velocity of the car is assumed to be known in the training dataset.

The third contribution is a mask on pixels that do not change between frames. These pixels are assumed to have no ego motion. Stationary pixels can occur for example if the car dashboard is visible in a frame, or other vehicles are moving at a similar speed nearby.

3.5 Digging Into Self-Supervised Monocular Depth Estimation [5]

The authors present MonoDepth2 on 17 Aug 2019. They propose mainly two contributions.

Firstly, instead of taking the average of the reprojection errors from all source frames given a target frame they use the minimum. This makes it so that if a feature is occluded in one source image but not in the other the errors will not be averaged together but instead the error from the source frame that is not occluded will be used.

Secondly, instead of calculating the loss for each depth scale in the decoder, all the depth maps are upsampled to the original target image size when computing the loss. This way a single pixel in the low resolution layer of the decoder will predict the depth of a patch of pixels in the originally sized input image.

3.6 Self-Supervised 3D Keypoint Learning for Ego-motion Estimation [12]

This paper from 7 Dec 2019 combines the work of the previously discussed papers and also adds keypoint learning from SuperPoint[1] and its successor Unsuperpoint[2] into the pipeline. The researchers train the depth, keypoint and pose estimating networks jointly, making them benefit from each other and achieve state of the art results.

4

Method

4.1 Unsupervised depth and ego motion learning

This section describes the methods used to learn how to predict depth and ego motion from a sequence of images.

4.1.1 Sequence datasets

The neural networks are trained and evaluated on 2 different datasets, KITTI[3] and Lyft[9]. Both datasets are preprocessed to remove frames where the camera is not moving. This is important because if there is no movement between frames then no depth information can be inferred during training when using monoscopic data. In the Kitti dataset the images from the left and right camera are treated as separate image sequences to yield more training data. The images are resized to 128x416 pixels, and the intrinsic camera matrix is updated accordingly. The lidar data from Kitti and Lyft is converted to sparse depth maps and are used during performance evaluation. The dataloader works by loading triplets of adjacent frames in the image sequences, and also the ground truth movement and depth map. An example of this is seen in Figure 4.1.



Figure 4.1: From top to bottom 3 frames from Kitti, I_{t-1} , I_t , I_{t+1} with the sparse depth map overlayed on frame I_t .

	Sequences / Samples	
	Train	Test
Kitti	110 / 16542	12 / 11349
Lyft	134 / 3759	14 / 1735

4.1.2 Architectures for depth and ego motion CNNs

In order to predict depth and motion from monocular images two different CNN architectures are examined.

The first architecture is referred to as SfMlearner[13]. The authors use a DispNet[10] architecture to predict depth maps at four different scales, and a ResNet18[7] architecture with modified decoder to predict pose updates in an euler angle axis representation.

The second architecture is referred to as Monodepth2[5]. The authors use a ResNet18 architecture instead of a DispNet architecture to predict depth estimates. They make this choice because its a smaller and faster architecture. Similarly they use a ResNet18 architecture with modified decoder to predict the pose updates in an euler angle axis representation.

4.1.3 Differentiable depth image warping

The core component of unsupervised depth learning is the differentiable depth image warp operation in the loss function of the CNN networks. Given the intrinsic camera matrix:

$$K = \begin{pmatrix} f_x & s & x_0 \\ 0 & f_y & y_0 \\ 0 & 0 & 1 \end{pmatrix}$$

And the predicted depth $D_t(p_t)$ of pixel p_t of the target (current) frame. And the transform $T_{t \rightarrow s}$ from the target to source (next/previous) frame:

$$T_{t \rightarrow s} = \begin{pmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 1 \end{pmatrix}$$

The position of the target pixel p_t in the source image p_s can be calculated in homogeneous coordinates as:

$$p_s \sim (K \quad \mathbf{0}) T_{t \rightarrow s} D_t(p_t) K^{-1} p_t$$

The pixel position p_s is however continuous and in order to sample the discrete source image I_s a differentiable bilinear sampling method is used. The method is described in *spatial transformer networks*[8] and works by interpolating the neighbouring 4 pixels values (top-left, bottom-right) by the distance to the the continuous sampling point p_t . This process is illustrated in Figure 4.2.

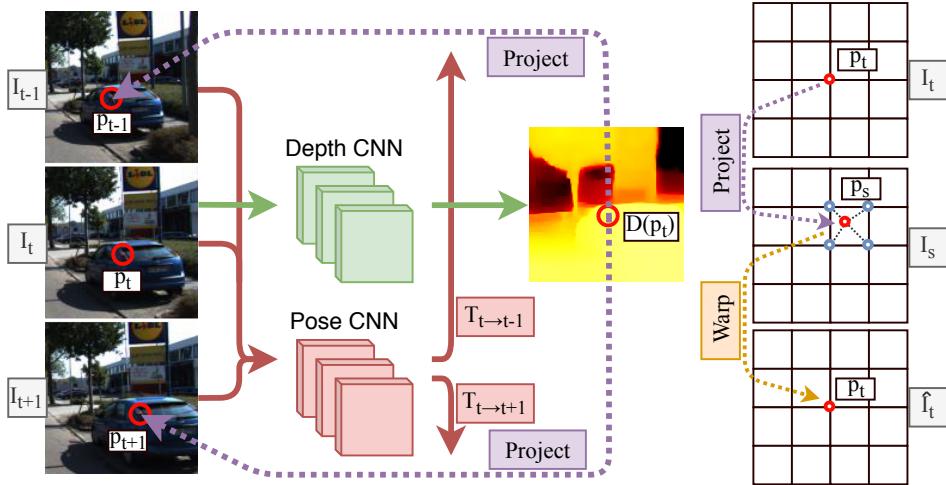


Figure 4.2: The CNN predicts the depth map D of the target image I_t , and also the relative movement, $T_{t \rightarrow t-1}$, $T_{t \rightarrow t+1}$ between the target image and the source images. Each pixel p_t in the target image is projected onto a position in the source images which are sampled using bilinear interpolation. This should recreate the appearance of the target image but with pixels sampled from the source image. An appearance similarity metric between the original target image and the recreated target images can be used as the loss function for the CNN to learn to accurately predict correct depth and movement.

4.1.4 Loss functions

In order to learn the objective of depth and ego motion prediction, different combinations of the following loss terms were evaluated.

To measure the similarity of the target image I_t and the reconstructed image \hat{I}_t a photometric loss term defined as the mean of the absolute value of the difference of pixel intensities of the two images was used.

$$\mathcal{L}_p(I_t, \hat{I}_t) = \text{mean}(|I_t - \hat{I}_t|)$$

Another photometric loss term evaluated in this thesis is based on structured similarity, referred to as SSIM[14]. It was originally developed to measure the quality of digital television, comparing a compressed digital image to the original distortion-free image.

$$\mathcal{L}_{ssim}(I_t, \hat{I}_t) = \text{clamp}(0, 1, \frac{1 - \text{SSIM}(I_t, \hat{I}_t)}{2})$$

$$\text{SSIM}(x, y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)}$$

with $C_1 = 1e-4$ and $C_2 = 9e-4$. To compute the per-patch mean μ_x and standard deviation σ_x I first used a 1 pixel reflection padding on the edges of the input images and then a 3×3 average pool filter with stride 1 to get the mean. Then $\sigma_x = \mu_{x^2} - \mu_x^2$.

The two above mentioned photometric loss terms can also be combined and balanced using $\mathcal{L}_{ps}(I_t, \hat{I}_s) = \alpha \mathcal{L}_{ssim} + (1 - \alpha) \mathcal{L}_p$. In my experiments I used $\alpha = 0.85$.

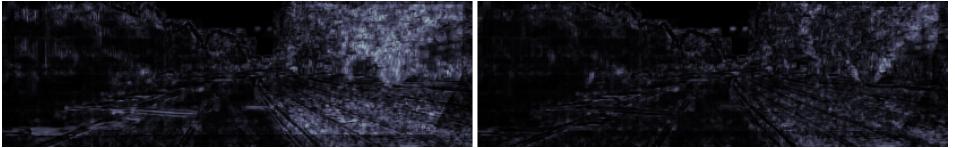


Figure 4.3: The images illustrates the photometric reconstruction loss \mathcal{L}_{ps} with $\alpha = 0.85$ for each pixel in the reconstructed image. The left image shows the loss after 1 epoch of training and the right image shows the loss after 30 epochs. The reconstruction loss should decrease during training as the network learns to predict better depth maps, which is what we see.

To propagate the depth from textured regions to regions of uniform color a depth smoothness loss term is used. The first alternative is a loss on the second derivative of the depth intensities.

$$\mathcal{L}_{smooth}(D_t) = \text{mean}(|\delta_x^2 D_t| + |\delta_y^2 D_t|)$$

I also ran experiments applying the smoothness loss on the disparity map instead of the depth map, but with worse performance. The second alternative investigated was an edge aware smoothness loss that weights the first order derivative of the depth map with the exponential of the first derivative of the pixel intensities. This method allows large changes in depth near sharp features in the image but penalizes changes in depth in smooth regions. The method showed great promise in the Monodepth2 paper, resulting in sharper edges because the smoothness term is weighted to mostly affect areas with small photometric derivative.

$$\mathcal{L}_{edge}(D_t) = \text{mean}(|\delta_x D_t| e^{-|\delta_x I_t|} + |\delta_y D_t| e^{-|\delta_y I_t|})$$

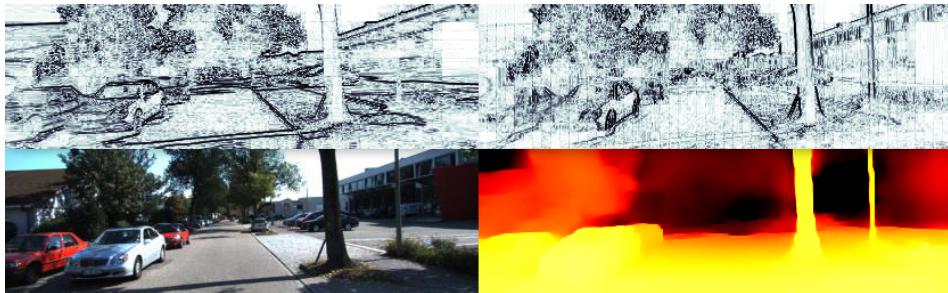


Figure 4.4: The top row of images illustrates the edge weighting terms $e^{-|\delta_y I_t|}$ and $e^{-|\delta_x I_t|}$ respectively. The weight is near 0 at the edges of tree trunks and near 1 on the pavement. This will preserve the details in the depth map around trees but keep the pavement smooth.

4.1.5 Handling occlusions

In the SfMLearner paper the photometric loss is calculated for the previous and next frames compared to the current in the sequence. The pixel wise average across the frames are then used. This causes problems if a pixel is for example occluded in the previous frame, but visible in the current and next frame. In this situation the average loss will be pretty high even though a correct depth and transformation has been predicted, because of the occluded pixel. Instead Monodepth2 suggests to pick the minimum per pixel error over the frames which alleviates this issue.

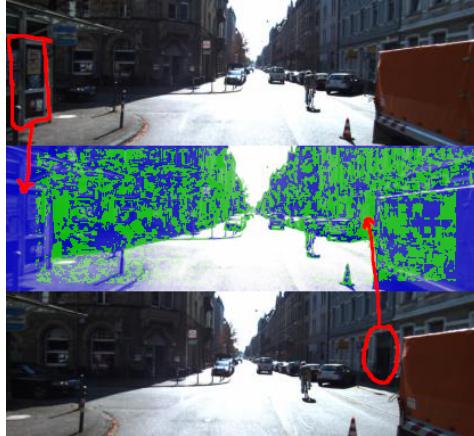


Figure 4.5: By picking the per pixel minimum reprojection error the issue created by occluded pixels can be alleviated. The top image is the previous frame, the middle image is the current target frame and the bottom image is the next frame in the sequence. If the minimum reprojection error can be found in the previous frame then it is colored blue, if its from the next frame it is green. Because the door to the right is occluded by the orange truck in the previous frame, the reprojection loss from the next frame is used instead where the door is visible. The wall to the left is outside the boundaries of the next image, so the reprojection error from the previous frame is used instead.

4.1.6 Handling model limitations

In order to optimize using the photometric reprojection error as the loss function two assumptions must hold. Firstly the scene must be static, meaning all objects in the scene must be still except the moving camera. Movement by cars and humans in the scene that is not due to the camera movement will cause problems. Secondly there must be photometric consistency between frames for the photometric error to make sense. This means that non lambertian surfaces, change in lighting, and change in exposure between frames will cause problems.

Explainability mask The authors of [13] tackle this problem by having a CNN predict what pixels are valid to use in the photometric loss function. It shares the encoder of the pose predicting network but branches off into a different encoder which estimates a mask of the valid/explainable pixels. The loss function for the mask is the cross entropy loss compared to a mask filled with ones. The photometric loss function is augmented to include the explainability mask removing pixels that cannot be explained by the predicted depth and transformation. This encourages the mask to be filled with ones, but allows some slack due to pixels that can not be explained by the photometric loss.

Stationary pixels mask The authors of [5] introduced a mask to remove stationary pixels from the set of previous, current and next frame. This is done by creating a mask where the photometric error is smaller before applying the projection than after. This works because stationary pixels that have not moved in relation to the camera will of course have a small photometric loss without reprojection. This will remove pixels from the car dashboard and also nearby vehicles that are traveling at the same speed.



Figure 4.6: The model assumes a stationary world without moving objects in it. By excluding stationary pixels from the loss issues arising from moving objects can be alleviated. The black pixels in the image are the ones removed because their photometric error is smaller before warping the image using the depth map compared to after warping. The red horizontal lines on the van and bicyclist illustrates that they do not move with respect to the camera, and the green slanted line illustrates that the back bike tire is moving with respect to the camera. The mask successfully removes pixels on the van and bicyclist, but also removes some pixels on the pavement that should not be removed.

4.2 Unsupervised keypoint learning

In order to build a map it is useful to extract features, or landmarks, from images. Every feature should have a unique descriptor which can be used to identify the features in the map in order to expand the map and localize the camera. How to build a map and localize in it is out of scope for this thesis. This section will describe an unsupervised learning method to extract features from images that would be usable in a full SfML pipeline.

The method evaluated in this thesis is based on the Unsuperpoint paper [2] from 2019. The network architecture is illustrated in Figure 4.7. The input image is

fed into a shared backbone. The features from the backbone are then split into three different encoders that estimate the score, position and a descriptor for each feature. The network will estimate a feature in every 8×8 patch of the image, but only the top N features sorted by score are used in the evaluation.

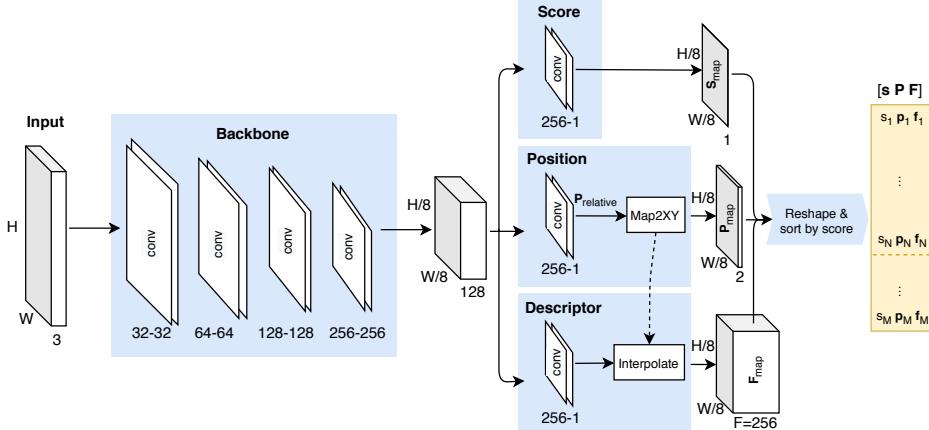


Figure 4.7: The network has a common backbone and then splits into separate score, position and descriptor encoders. The output is reshaped and sorted by descending score.

The network is trained in a siamese twin setup, where two duplicate networks that share weights are fed different inputs and the loss is formulated by comparing the output scores, positions and descriptors.

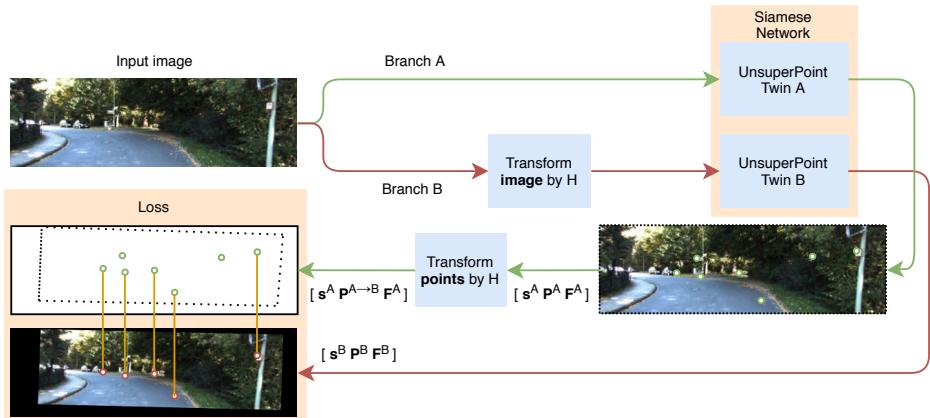


Figure 4.8

The backbone consists of four groups of convolution, batch norm, ReLu and max

pooling operations.

4.2.1 Score encoder

The score encoder

Brute force matching using PyTorch instead of OpenCV because we want to use the result as input to the next consensus step.

4.3 Unsupervised consensus maximization

Consensus maximization in a set of features related by a geometric constraint that also contains outliers is an important aspect of an SFM system. The problem consists of finding the largest set of inlier features in a set that is consistent with a geometric constraint. In a 3D vision setting the geometric constraint could be for example a rigid transformation $[R \ t]$, a fundamental matrix F or a homography H that relates pairs of corresponding 3D or 2D points.

A popular algorithm to solve this problem is RANSAC[?], but is however not differentiable and therefore not learnable by a neural network. In this paper[?] from 2019 the authors present a differential method to learn consensus maximization in an unsupervised setting. The explanation of the method in the paper is however very theoretical and non-accessible to master level students. In this section I will try to explain the method with emphasis on ease of understanding and with the mathematical toolbox of an engineering student. I will focus on the method used to extract the parameters of the constraint, as it's not explained by the paper in detail and took a lot of figuring out on my part, especially how to extract the homography.

Consider one of the simplest geometric constraint on a set of corresponding 2D points in subsequent frames of the KITTI dataset, the fundamental matrix F . Given a point u in frame 1, and point v in frame 2 they are related in the following way.

$$u^T F v = 0$$

To be more explicit the relation can be expressed as follows.

$$\begin{pmatrix} u_x & u_y & 1 \end{pmatrix} \begin{pmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{pmatrix} \begin{pmatrix} v_x \\ v_y \\ 1 \end{pmatrix} = 0$$

Multiplication of the two rightmost matrices gives:

$$\begin{pmatrix} u_x & u_y & 1 \end{pmatrix} \begin{pmatrix} f_{11}v_x + f_{12}v_y + f_{13} \\ f_{21}v_x + f_{22}v_y + f_{23} \\ f_{31}v_x + f_{32}v_y + f_{33} \end{pmatrix} = 0$$

Multiplication of the remaining two matrices and grouping the monomials in parentheses for readability gives the following linear equation.

$$f_{11}(u_x v_x) + f_{12}(u_x v_y) + f_{13}(u_x) + f_{21}(u_y v_x) + f_{22}(u_y v_y) + f_{23}(u_y) + f_{31}(v_x) + f_{32}(v_y) + f_{33}(1) = 0$$

Now consider a set of m corresponding points that are related by the same fundamental matrix F . Place the 9 monomials in the parentheses for all m points in a matrix $M \in \mathbb{R}^{m \times 9}$. A matrix of monomials, such as this one, is called a Vandermonde matrix.

$$M = \begin{pmatrix} u_{x,1}v_{x,1} & u_{x,1}v_{y,1} & u_{x,1} & u_{y,1}v_{x,1} & u_{y,1}v_{y,1} & u_{y,1} & v_{x,1} & v_{y,1} & 1 \\ & & & & \vdots & & & & \\ u_{x,m}v_{x,m} & u_{x,m}v_{y,m} & u_{x,m} & u_{y,m}v_{x,m} & u_{y,m}v_{y,m} & u_{y,m} & v_{x,m} & v_{y,m} & 1 \end{pmatrix}$$

The constraint is satisfied by all m points if

$$M \begin{pmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}^{mx1}$$

Now the solution \mathbf{f} should look very familiar. The column vector \mathbf{f} is in fact the very definition of a nullspace of M . This is the key insight in this method of unsupervised consensus maximization. Finding the nullspace of M is done by performing a singular value decomposition, which is a differentiable operation.

$$U\Sigma V^T = M$$

The nullspace will be the last column vector of $V \in \mathbb{R}^{9 \times 9}$ and the last singular value σ_9 in the diagonal of Σ will be zero. This is correct only if the set of corresponding points contains only inliers. If the set contains outliers that can not be

related by the geometric constraint then the last singular value σ_9 will be greater than zero. If the inliers and outliers were known they could be represented by a vector $\mathbf{w} \in [0, 1]^m$ where 1 means inlier and 0 means outlier. Now M can be weighted by \mathbf{w} which would cancel all rows in M which comes from an outlier pair of points.

$$\text{diag}(\mathbf{w})M\mathbf{f} = \mathbf{0}$$

The problem of finding the largest set of inliers has now been transformed into finding the weights \mathbf{w} such that the last singular value σ_9 of $\text{diag}(\mathbf{w})M$ is near zero, and the sum of the weights should also be large to have as many inliers as possible.

To learn the weight $w \in \mathbf{w}$ for each pair of points, the PointNet[?] architecture was used. The pair of points are concatenated into rows of a $X \in \mathbb{R}^{mxn}$ matrix. For 2D points $n = 2 + 2 = 4$ and for 3D points $n = 3 + 3 = 6$. The input X is fed through the PointNet segmentation network with a single output class, signifying if the point pair is an inlier or not. The key innovation of the PointNet architecture is that it can directly consume point cloud data, as opposed to first transforming the input to a 3D voxel grid or an image. The network is also invariant to the ordering of the points in the input list. Figure 4.9 illustrates the process of feeding the input through pointnet to get the weights, and then applying the weights to the Vandermonde matrix of which we want to minimize the last singular value.

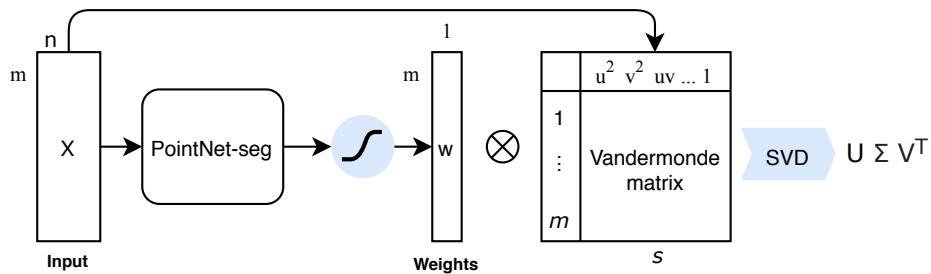


Figure 4.9

Because the fundamental matrix is a constraint expressed by 1 linear equation a basis of only 1 nullspace vector is needed from the singular value decomposition to construct it. But the method can also be used to learn constraints such as a homography or rigid 3D transformation that are constrained by 3 linear equations and therefore require a basis of 3 nullspace vectors. In that case the last 3 singular values from the SVD should be minimized in the loss function.

In the following sections the method used to extract the homography from the basis vectors will be explained. Two points that are related by a homography is

expressed as follows. The monomials are grouped by parentheses.

$$\begin{aligned}
 u \times Hv = 0 &\rightarrow \\
 [u]_x Hv = 0 &\rightarrow \\
 \begin{pmatrix} 0 & -1 & u_y \\ 1 & 0 & -u_x \\ -u_y & u_x & 0 \end{pmatrix} \begin{pmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{pmatrix} \begin{pmatrix} v_x \\ v_y \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \rightarrow \\
 \begin{cases} -h_{21}(v_x) - h_{22}(v_y) - h_{23}(1) + h_{31}(v_x u_y) + h_{32}(v_y u_y) + h_{33}(u_y) = 0 \\ h_{11}(v_x) + h_{12}(v_y) + h_{13}(1) - h_{31}(v_x u_x) - h_{32}(v_y u_x) - h_{33}(u_x) = 0 \\ -h_{11}(v_x u_y) - h_{12}(v_y u_y) - h_{13}(u_y) + h_{21}(v_x u_x) + h_{22}(v_y u_x) + h_{23}(u_x) = 0 \end{cases}
 \end{aligned}$$

The monomials are the same as for the fundamental constraint, which means that the same Vandermonde matrix M can be used. But because an homography is constrained by $r = 3$ linear equations the last 3 singular values of the SVD should be minimized during training.

Extract 3 basis vectors from the nullspace which will be the 3 rightmost columns of V in the SVD of $\text{diag}(\mathbf{w})M$.

$$B = (\mathbf{v}_7 \quad \mathbf{v}_8 \quad \mathbf{v}_9) \in \mathbb{R}^{9 \times 3}$$

The 3 linear equations of the homography are tangled which means that the elements of B can not be directly mapped onto the elements in H . Notice that the first equation of H does not contain u_y and the second equation of H does not contain u_x . We can exploit this fact to perform a change of basis from $B \in \mathbb{R}^{9 \times 3}$ to $B' \in \mathbb{R}^{9 \times 2}$ that should have the following structure.

Monomial	Structure of B'	Corresponding element in H
$u_x v_x$	$\begin{pmatrix} 0 & . \\ 0 & . \\ 0 & . \end{pmatrix}$	$\begin{pmatrix} h_{31} \\ h_{32} \\ h_{33} \end{pmatrix}$
$u_x v_y$	$\begin{pmatrix} . & 0 \\ . & 0 \\ . & 0 \end{pmatrix}$	$\begin{pmatrix} h_{31} \\ h_{32} \\ h_{33} \end{pmatrix}$
u_x	$\begin{pmatrix} . & . \\ . & . \\ . & . \end{pmatrix}$	$\begin{pmatrix} h_{21} & h_{11} \\ h_{22} & h_{12} \\ h_{23} & h_{13} \end{pmatrix}$
$u_y v_x$	$\begin{pmatrix} . & 0 \\ . & 0 \\ . & 0 \end{pmatrix}$	$\begin{pmatrix} h_{31} \\ h_{32} \\ h_{33} \end{pmatrix}$
$u_y v_y$	$\begin{pmatrix} . & . \\ . & . \\ . & . \end{pmatrix}$	$\begin{pmatrix} h_{31} \\ h_{32} \\ h_{33} \end{pmatrix}$
u_y	$\begin{pmatrix} . & . \\ . & . \\ . & . \end{pmatrix}$	$\begin{pmatrix} h_{21} & h_{11} \\ h_{22} & h_{12} \\ h_{23} & h_{13} \end{pmatrix}$
v_x	$\begin{pmatrix} . & . \\ . & . \\ . & . \end{pmatrix}$	$\begin{pmatrix} h_{21} & h_{11} \\ h_{22} & h_{12} \\ h_{23} & h_{13} \end{pmatrix}$
v_y	$\begin{pmatrix} . & . \\ . & . \\ . & . \end{pmatrix}$	$\begin{pmatrix} h_{21} & h_{11} \\ h_{22} & h_{12} \\ h_{23} & h_{13} \end{pmatrix}$
1	$\begin{pmatrix} . & . \\ . & . \\ . & . \end{pmatrix}$	$\begin{pmatrix} h_{21} & h_{11} \\ h_{22} & h_{12} \\ h_{23} & h_{13} \end{pmatrix}$

To perform the change of basis to get the structure of B' we use the nullspace $\mathbf{n}_1 \in \mathbb{R}^{3 \times 1}$ of rows 1, 2 and 3 from B , and the nullspace $\mathbf{n}_2 \in \mathbb{R}^{3 \times 1}$ of row 4, 5, 6 from B .

$$A_1 = \begin{pmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{pmatrix}, A_2 = \begin{pmatrix} b_{41} & b_{42} & b_{43} \\ b_{51} & b_{52} & b_{53} \\ b_{61} & b_{62} & b_{63} \end{pmatrix}$$

\mathbf{n}_1 = "rightmost column of right-singular vectors of A1"

\mathbf{n}_2 = "rightmost column of right-singular vectors of A2"

$$\mathbf{b}^{n_1} = B\mathbf{n}_1$$

$$\mathbf{b}^{n_2} = B\mathbf{n}_2$$

Now the \mathbf{b}^{n_1} and \mathbf{b}^{n_2} basis vectors will have the following structure as desired.

$$\mathbf{b}^{n_1} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{pmatrix}, \mathbf{b}^{n_2} = \begin{pmatrix} \cdot \\ \cdot \\ \cdot \\ 0 \\ 0 \\ 0 \\ \cdot \end{pmatrix},$$

The basis vectors have zeros at the correct place. The next step is to adjust the scale so that row 4, 5 and 6 of \mathbf{b}^{n_1} and row 1, 2 and 3 of \mathbf{b}^{n_2} have the same norm because they both represent the same elements h_{31} , h_{32} and h_{33} . We also make sure they have the same sign.

$$s = \text{sign}(b_1^{n_2} + b_2^{n_2} + b_3^{n_2}) \text{sign}(b_4^{n_1} + b_5^{n_1} + b_6^{n_1})$$

s will be -1 if they have different signs, or 1 if they are the same sign.

$$B' = \left(\mathbf{b}^{n_1} / \| (b_4^{n_1} \ b_5^{n_1} \ b_6^{n_1}) \| \quad \mathbf{b}^{n_2} / \| (b_1^{n_2} \ b_2^{n_2} \ b_3^{n_2}) \| \ s \right) \in \mathbb{R}^{9 \times 2}$$

It is now possible to assign the elements of H by pattern matching.

$$H = \begin{pmatrix} b'_{72} & b'_{82} & b'_{92} \\ b'_{71} & b'_{81} & b'_{91} \\ -b'_{41} & -b'_{51} & -b'_{61} \end{pmatrix}$$

4.4 Evaluation

This section describes the metrics used to evaluate depth, ego motion, feature point and consensus maximization predictions. These specific metrics were chosen because they are also used in the related works in this field, which makes the results comparable to other papers.

4.4.1 Depth

The depth error and accuracy metrics are sparse and only calculated for the pixels of which there exists a ground truth laser measurement in the dataset. The values are averaged over all laser measurements and frames in the test split of the dataset. An error of 0 and an accuracy of 1 is optimal, but can never be achieved in practice.

The depth is predicted relative to an unknown scale, and the ground truth is measured in meters. To alleviate this issue the depth predictions are scaled so that their median is the same as the median of the ground truth for each frame. This does not guarantee that the unit of the predictions is meters, but at least it makes the scales somewhat similar.

The error metrics used in the results chapter are referred to as *Abs Rel*, *Sq Rel*, *RMSE* and *RMSLE*.

The *Abs Rel* error is based on *MAE* which is the mean of the absolute errors which means that it has the same unit as the errors, and is conceptually quite easy to interpret.

$$\text{MAE} = \frac{\sum_{n=1}^N |y_n - \hat{y}_n|}{N}$$

Because the depth from the neural network is without unit and only predicted relative to an unknown scale, a variation on the *MSE* metric called *Abs Rel* is used.

$$\text{Abs Rel} = \frac{\sum_{n=1}^N \frac{|y_n - \hat{y}_n|}{y_n}}{N}$$

The *MSE* metric is the mean of squared errors, for an unbiased estimator it represents the variance of the errors.

$$\text{MSE} = \frac{\sum_{n=1}^N (y_n - \hat{y}_n)^2}{N}$$

Once again, because the depth is predicted relative to an unknown scale a variation on *MSE* called *Sq Rel* is used instead.

$$\text{Sq Rel} = \frac{\sum_{n=1}^N \frac{(y_n - \hat{y}_n)^2}{y_n}}{N}$$

The *RMSE* metric is the square root of the mean of squared errors. For an unbiased estimator it represents the standard deviation of the errors. Because the

errors are squared in the *RMSE* metric it is more sensitive to outliers than for example *MSE*.

$$\text{RMSE} = \sqrt{\frac{\sum_{n=1}^N (y_n - \hat{y}_n)^2}{N}}$$

The *RMSLE* metric is similar to *RMSE* but is useful because it penalizes large errors less when both the actual and predicted values are large.

$$\text{RMSLE} = \sqrt{\frac{\sum_{n=1}^N (\log y_n - \log \hat{y}_n)^2}{N}}$$

To measure the depth accuracy, in the range from 0 to 1, the following metric is used.

$$a_\gamma = \frac{\sum_{n=1}^N (\max(\frac{y_n}{\hat{y}_n}, \frac{\hat{y}_n}{y_n}) < \gamma)}{N}, \text{ for } \gamma \in \{1.25, 1.25^2, 1.25^3\}$$

This should be interpreted as the ratio of predictions that is within the ratio of γ relative to the ground truth.

4.4.2 Ego motion

The error of the camera motion predictions are measured using *RMSE*. But instead of taking the mean over all poses in a sequence, the alignment error is calculated part wise over a track length of only 5 poses. The final error is presented as the mean *RMSE* of all parts. Each track part is transformed to have its first pose coincide with the world origo. This makes it so that the first pose in the track part of both the prediction and ground truth is the identity matrix $I \in \mathbb{R}^{4 \times 4}$. Because the ego motion is predicted relative to an unknown scale and the ground truth is in meters, we scale each predicted track part to have similar scale to the ground truth track part.

4.4.3 Keypoints

4.4.4 Consensus maximization

5

Results

Figure 5.2 contains some preliminary evaluation results. The gray configurations in Figure 5.1 have not been trained yet. Each configuration takes about 15 hours to train. Still missing possibility to train on Synthia and also random data augmentations, L1 loss on disparity and velocity supervision loss.

Name	Architecture	Training Dataset	Smoothing			Auto masking		SSIM	Avg/Min	Upscale
			Which	Edge aware	Normalize	Predict	Stationary			
B0	sfmlearner	kitti	disp			x			avg	
B1	sfmlearner	kitti	depth			x			avg	
B2	sfmlearner	kitti	depth						avg	
B3	sfmlearner	kitti	depth				x		avg	
B4	sfmlearner	kitti	depth	x			x		avg	
B5	sfmlearner	kitti	depth	x	x	x	x		min	
B6	sfmlearner	kitti	depth	x	x	x	x		min	
B7	sfmlearner	kitti	depth	x	x	x	x		min	x
B8	sfmlearner	lyft	depth	x	x	x	x		min	x
C3	monodepth2	kitti	depth			x			avg	
C6	monodepth2	kitti	depth	x	x	x	x		min	
C7	monodepth2	kitti	depth	x	x	x	x		min	x
C8	monodepth2	lyft	depth	x	x	x	x		min	x

Figure 5.1: Different configurations of network architectures, training datasets and loss terms evaluated to find the best performance

Configuration	Test Dataset	Depth error metric				Depth accuracy metric			Trajectory error metric RMSE
		Abs Rel	Sq Rel	RMSE	RMSE log	a < 1.25	a < 1.25^2	a < 1.25^3	
B0	kitti	0,362	13,469	8,759	0,372	0,73	0,884	0,933	0,02
B1	kitti	0,184	1,718	4,953	0,258	0,776	0,915	0,96	0,02
B2	kitti	0,174	1,405	4,829	0,249	0,784	0,92	0,964	0,024
B3	kitti								
B4	kitti								
B5	kitti								
B6	kitti	0,137	0,797	4,282	0,208	0,837	0,948	0,977	0,024
B7	kitti								
B7	lyft								
B8	kitti								
B8	lyft								
C3	kitti								
C6	kitti	0,126	0,714	4,018	0,194	0,86	0,958	0,982	0,019
C7	kitti	0,132	0,769	3,966	0,196	0,859	0,957	0,981	0,019
C7	lyft	0,304	7,019	21,907	0,414	0,518	0,775	0,886	0,042
C8	kitti								
C8	lyft								
SfMLearner	kitti	0,208	1,768	6,856	0,283	0,678	0,885	0,957	
Monodepth2	kitti	0,132	1,044	5,142	0,21	0,845	0,948	0,977	

Figure 5.2: Evaluation metrics when testing the configurations on the testing split of the datasets

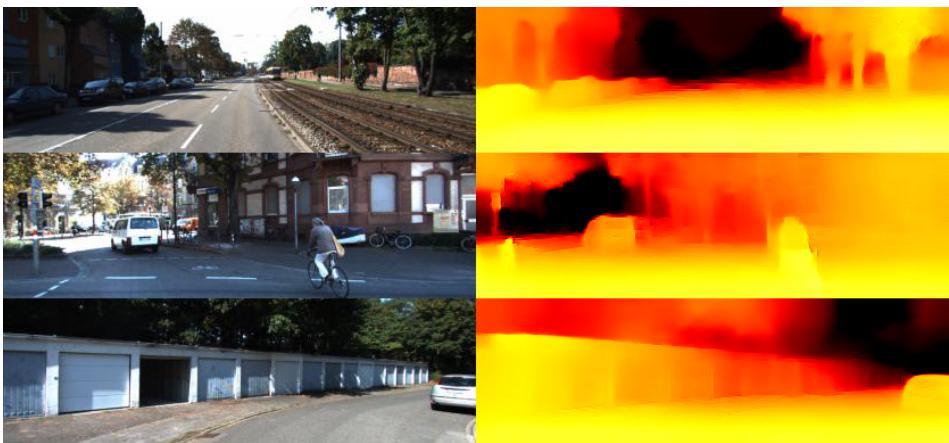


Figure 5.3: Examples from the Kitti dataset



Figure 5.4: 3D render of colorized depth map

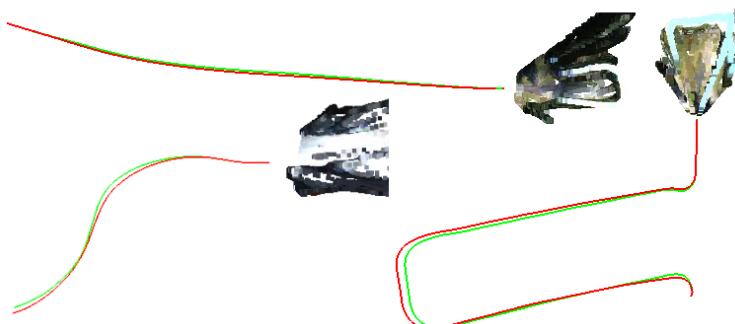


Figure 5.5: 3D visualization of the camera movement in three different image sequences. The green lines are the ground truth and the red lines are the predicted camera trajectories.



Figure 5.6

6

Discussion

The best combination for depth and ego motion predictions.

Small disparity, large depth, for cars moving with the same relative speed, no ego motion.

TODO: Image of holes from small disparity

Over-fitting while training homography.

TODO: Plot over fitting. Image + Loss + Histogram

Alterations needed to make consensus converge.

Future work to couple the training of depth, keypoints and consensus more tightly.

Bibliography

- [1] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. 12 2017.
- [2] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. pages 337–33712, 06 2018. doi: 10.1109/CVPRW.2018.00060.
- [3] Andreas Geiger, P Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: the kitti dataset. *The International Journal of Robotics Research*, 32:1231–1237, 09 2013. doi: 10.1177/0278364913491297.
- [4] Clement Godard, Oisin Aodha, and Gabriel Brostow. Unsupervised monocular depth estimation with left-right consistency. 07 2017. doi: 10.1109/CVPR.2017.699.
- [5] Clément Godard, Oisin Aodha, and Gabriel Brostow. Digging into self-supervised monocular depth estimation, 06 2018.
- [6] Vitor Guizilini, Rares Ambrus, Sudeep Pillai, and Adrien Gaidon. Packnet-sfm: 3d packing for self-supervised monocular depth estimation, 05 2019.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. 7, 12 2015.
- [8] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. *Advances in Neural Information Processing Systems 28 (NIPS 2015)*, 06 2015.
- [9] R. Kesten, M. Usman, J. Houston, T. Pandya, K. Nadhamuni, A. Ferreira, M. Yuan, B. Low, A. Jain, P. Ondruska, S. Omari, S. Shah, A. Kulkarni, A. Kazakova, C. Tao, L. Platinsky, W. Jiang, and V. Shet. Lyft level 5 av dataset 2019. url`https://level5.lyft.com/dataset/`, 2019.
- [10] N.Mayer, E.Ilg, P.Häusser, P.Fischer, D.Cremers, A.Dosovitskiy, and T.Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *IEEE International*

- Conference on Computer Vision and Pattern Recognition (CVPR), 2016.*
URL <http://lmb.informatik.uni-freiburg.de/Publications/2016/MIFDB16>. arXiv:1512.02134.
- [11] Sudeep Pillai, Rares Ambrus, and Adrien Gaidon. Superdepth: Self-supervised, super-resolved monocular depth estimation. pages 9250–9256, 05 2019. doi: 10.1109/ICRA.2019.8793621.
 - [12] Jie Xiong Tang, Rares Ambrus, Vitor Guizilini, Sudeep Pillai, Hanme Kim, and Adrien Gaidon. Self-supervised 3d keypoint learning for ego-motion estimation, 12 2019.
 - [13] Tinghui Zhou, Matthew Brown, Noah Snavely, and David Lowe. Unsupervised learning of depth and ego-motion from video. 04 2017.
 - [14] Zhou Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.